

Pilastro dell'efficienza delle prestazioni



Pilastro dell'efficienza delle prestazioni: Framework AWS Well-Architected

Copyright © 2025 Amazon Web Services, Inc. and/or its affiliates. All rights reserved.

I marchi e l'immagine commerciale di Amazon non possono essere utilizzati in relazione a prodotti o servizi che non siano di Amazon, in una qualsiasi modalità che possa causare confusione tra i clienti o in una qualsiasi modalità che denigri o discrediti Amazon. Tutti gli altri marchi non di proprietà di Amazon sono di proprietà delle rispettive aziende, che possono o meno essere associate, collegate o sponsorizzate da Amazon.

Table of Contents

Riassunto e introduzione	1
Introduzione	1
Efficienza delle prestazioni	3
Principi di progettazione	3
Definizione	4
Scelta dell'architettura	5
PERF01-BP01 Scopri e comprendi i servizi e le funzionalità cloud disponibili	5
Guida all'implementazione	6
Risorse	7
PERF01-BP02 Utilizzo delle indicazioni del provider cloud o di un partner appropriato per conoscere gli schemi di architettura e le best practice	8
Guida all'implementazione	6
Risorse	7
PERF01-BP03 Fattore di costo nelle decisioni architettoniche	10
Guida all'implementazione	6
Risorse	7
PERF01-BP04 Valuta l'impatto dei compromessi sui clienti e sull'efficienza dell'architettura	12
Guida all'implementazione	6
Risorse	7
PERF01-BP05 Utilizza politiche e architetture di riferimento	14
Guida all'implementazione	6
Risorse	7
PERF01-BP06 Uso del benchmarking per guidare le decisioni sull'architettura	15
Guida all'implementazione	6
Risorse	7
PERF01-BP07 Utilizza un approccio basato sui dati per le scelte architettoniche	18
Guida all'implementazione	6
Risorse	7
Calcolo e hardware	21
PERF02-BP01 Seleziona le migliori opzioni di elaborazione per il tuo carico di lavoro	21
Guida all'implementazione	6
Passaggi dell'implementazione	6
Risorse	7
PERF02-BP02 Comprendi la configurazione e le funzionalità di elaborazione disponibili	25

Guida all'implementazione	6
Passaggi dell'implementazione	6
Risorse	7
PERF02-BP03 Raccogli metriche relative al calcolo	29
Guida all'implementazione	6
Passaggi dell'implementazione	6
Risorse	7
PERF02-BP04 Configurazione e dimensionamento corretto delle risorse di calcolo	31
Guida all'implementazione	6
Risorse	7
PERF02-BP05 Scala le tue risorse di elaborazione in modo dinamico	34
Guida all'implementazione	6
Risorse	7
PERF02-BP06 Uso di acceleratori di elaborazione ottimizzati basati su hardware	37
Guida all'implementazione	6
Risorse	7
Gestione dei dati	40
PERF03-BP01 Utilizza un archivio dati appositamente progettato che supporti al meglio i requisiti di accesso e archiviazione dei dati	40
Guida all'implementazione	6
Risorse	7
PERF03-BP02 Valuta le opzioni di configurazione disponibili per l'archivio dati	52
Guida all'implementazione	6
Risorse	7
PERF03-BP03 Raccogli e registra le metriche delle prestazioni degli archivi dati	57
Guida all'implementazione	6
Passaggi dell'implementazione	6
Risorse	7
PERF03-BP04 Implementazione di strategie per migliorare le prestazioni delle query nel datastore	60
Guida all'implementazione	6
Risorse	7
PERF03-BP05 Implementa modelli di accesso ai dati che utilizzano la memorizzazione nella cache	62
Guida all'implementazione	6
Risorse	7

Reti e distribuzione di contenuti	66
PERF04-BP01 Scopri come la rete influisce sulle prestazioni	66
Guida all'implementazione	6
Risorse	7
PERF04-BP02 Valuta le funzionalità di rete disponibili	70
Guida all'implementazione	6
Risorse	7
PERF04-BP03 Scegli la connettività dedicata appropriata o per il tuo carico di lavoro VPN	77
Guida all'implementazione	6
Risorse	7
PERF04-BP04 Utilizza il bilanciamento del carico per distribuire il traffico su più risorse	80
Guida all'implementazione	6
Risorse	7
PERF04-BP05 Scegli i protocolli di rete per migliorare le prestazioni	84
Guida all'implementazione	6
Risorse	7
PERF04-BP06 Scegli la posizione del carico di lavoro in base ai requisiti di rete	88
Guida all'implementazione	6
Risorse	7
PERF04-BP07 Ottimizza la configurazione di rete in base a metriche	93
Guida all'implementazione	6
Risorse	7
Processo e cultura	98
PERF05-BP01 Stabilisci indicatori chiave di prestazione () per misurare lo stato e le prestazioni del carico di lavoro KPIs	100
Guida all'implementazione	6
Passaggi dell'implementazione	6
Risorse	7
PERF05-BP02 Utilizza soluzioni di monitoraggio per comprendere le aree in cui le prestazioni sono più critiche	103
Guida all'implementazione	6
Risorse	7
PERF05-BP03 Definire un processo per migliorare le prestazioni del carico di lavoro	106
Guida all'implementazione	6
Risorse	7
PERF05-BP04 Load Esegui un test del tuo carico di lavoro	107

Guida all'implementazione	6
Risorse	7
PERF05-BP05 Usa l'automazione per risolvere in modo proattivo i problemi relativi alle prestazioni	110
Guida all'implementazione	6
Risorse	7
PERF05-BP06 Conserva il carico di lavoro e i servizi up-to-date	112
Guida all'implementazione	6
Passaggi dell'implementazione	6
Risorse	7
PERF05-BP07 Rivedi le metriche a intervalli regolari	114
Guida all'implementazione	6
Risorse	7
Conclusioni	117
Collaboratori	118
Approfondimenti	119
Revisioni del documento	120
Note	122
AWS Glossario	123

Pilastro dell'efficienza delle prestazioni: Framework AWS Well-Architected

Data di pubblicazione: 6 novembre 2024 ([Revisioni del documento](#))

Il presente whitepaper riguarda il pilastro dell'efficienza delle prestazioni del Framework AWS Well-Architected. Inoltre, fornisce indicazioni per aiutarti i clienti ad applicare le best practice per la progettazione, la distribuzione e la manutenzione degli AWS ambienti.

Introduzione

Il [Framework AWS Well-Architected](#) aiuta a comprendere i pro e i contro delle decisioni prese durante la creazione dei carichi di lavoro in AWS. Utilizzando il Framework, scoprirai le best practice architetturali per progettare e gestire carichi di lavoro affidabili, sicuri, efficienti, convenienti e sostenibili nel cloud. Il Framework permette di misurare in modo coerente le architetture secondo le best practice e identificare le aree di miglioramento. Disporre di carichi di lavoro ben progettati aumenta notevolmente la probabilità di successo aziendale.

Il Framework si basa su sei pilastri:

- Eccellenza operativa
- Sicurezza
- Affidabilità
- Efficienza delle prestazioni
- Ottimizzazione dei costi
- Sostenibilità

Il presente whitepaper tratta dell'applicazione dei principi del pilastro dell'efficienza delle prestazioni ai carichi di lavoro. Nei tradizionali ambienti on-premises, raggiungere prestazioni durature e di alto livello può essere difficoltoso. L'utilizzo dei principi contenuti in questo documento ti aiuterà a creare architetture in AWS in grado di offrire prestazioni efficienti e sostenibili nel tempo. Linee guida e best practice contenute nel presente documento sono suddivise in cinque aree di interesse chiave, che costituiscono principi guida per la creazione di soluzioni cloud in AWS efficienti in termini di prestazioni. Queste aree di interesse sono:

- [Scelta dell'architettura](#)
- [Calcolo e hardware](#)
- [Gestione dei dati](#)
- [Reti e distribuzione di contenuti](#)
- [Processo e cultura](#)

Questo documento è rivolto ai ruoli nell'ambito della tecnologia, ad esempio a Chief Technology Officer (CTO), progettisti, sviluppatori e membri dei team operativi. Dopo avere letto questo documento, comprenderai le best practice di AWS e le strategie da utilizzare durante la progettazione di architetture di un ambiente cloud ad alte prestazioni.

Efficienza delle prestazioni

Il pilastro dell'efficienza delle prestazioni include la capacità di utilizzare in modo efficiente le risorse nel cloud per soddisfare i requisiti in termini di prestazione e di mantenere tale efficienza a fronte al cambiamento della domanda e all'evoluzione delle tecnologie.

Argomenti

- [Principi di progettazione](#)
- [Definizione](#)

Principi di progettazione

I seguenti principi di progettazione possono aiutarti a raggiungere e mantenere carichi di lavoro efficienti nel cloud.

- **Estendi a tutti le tecnologie avanzate:** facilita l'implementazione di tecnologie avanzate da parte del tuo team delegando le attività complesse al tuo fornitore di cloud. Anziché chiedere al team IT di imparare come adottare e gestire una nuova tecnologia, valuta l'opportunità di utilizzare la tecnologia come servizio. Ad esempio, No SQL database, transcodifica multimediale e apprendimento automatico sono tutte tecnologie che richiedono competenze specialistiche. Nel cloud, tali tecnologie diventano servizi che il tuo team può semplicemente utilizzare mentre si concentra sullo sviluppo di un prodotto invece che sul provisioning e sulla gestione delle risorse.
- **Diventa globale in pochi minuti:** l'implementazione del carico di lavoro in più AWS regioni del mondo ti consente di fornire una latenza inferiore e un'esperienza migliore per i tuoi clienti a costi minimi.
- **Utilizza architetture serverless:** scegliendo le architetture serverless, non avrai più bisogno di gestire e mantenere server fisici per portare a termine le attività di elaborazione tradizionali. Ad esempio, i servizi di storage serverless possono agire da siti web statici, eliminando la necessità di server web, mentre i servizi di eventi possono ospitare il codice. Questo elimina l'onere operativo della gestione dei server fisici, con una riduzione dei costi delle transazioni, dal momento che servizi gestiti di questo tipo funzionano a livello di cloud.
- **Sperimenta più di frequente:** le risorse virtuali e automatizzabili ti permettono di portare a termine velocemente i test comparativi utilizzando diversi tipi di istanze, storage o configurazioni.

- Prendi in considerazione la comprensione meccanica: sfrutta la strategia tecnologica più adatta ai tuoi obiettivi. Ad esempio, prendi in considerazione gli schemi di accesso ai dati quando scegli una strategia basata su database o archiviazione per il tuo carico di lavoro.

Definizione

Concentrati sulle seguenti aree per ottenere l'efficienza delle prestazioni nel cloud:

- [Scelta dell'architettura](#)
- [Calcolo e hardware](#)
- [Gestione dei dati](#)
- [Reti e distribuzione di contenuti](#)
- [Processo e cultura](#)

Adotta un approccio basato sui dati per creare un'architettura ad alte prestazioni. Raccogli dati su tutti gli aspetti dell'architettura, dalla progettazione di alto livello alla selezione e alla configurazione dei tipi di risorse.

Rivedendo regolarmente le tue scelte, ti assicurerai di sfruttare i vantaggi del cloud in continua evoluzione. AWS Il monitoraggio ti assicurerà di essere consapevole di qualsiasi divergenza rispetto alle prestazioni previste. Infine, puoi raggiungere dei compromessi nella tua architettura per migliorare le prestazioni, per esempio utilizzando la compressione o la memorizzazione nella cache oppure allentando i requisiti di coerenza.

Scelta dell'architettura

La soluzione ottimale per un determinato carico di lavoro può variare e le soluzioni spesso combinano molteplici approcci. I carichi di lavoro Well-Architected utilizzano soluzioni multiple e forniscono funzionalità diverse per migliorare le prestazioni.

Le risorse AWS sono disponibili in diverse configurazioni e tipologie, il che semplifica la ricerca di un approccio che soddisfi appieno le tue esigenze. Inoltre, puoi trovare opzioni che non sono facili da trovare nelle infrastrutture on-premises. Ad esempio, un servizio gestito come Amazon DynamoDB offre un database NoSQL interamente gestito, con una latenza di pochissimi millisecondi, indipendentemente dalle dimensioni.

Questa area di interesse offre linee guida e best practice su come selezionare risorse cloud e modelli di architettura efficienti e ad alte prestazioni.

Best practice

- [PERF01-BP01 Scopri e comprendi i servizi e le funzionalità cloud disponibili](#)
- [PERF01-BP02 Utilizzo delle indicazioni del provider cloud o di un partner appropriato per conoscere gli schemi di architettura e le best practice](#)
- [PERF01-BP03 Fattore di costo nelle decisioni architettoniche](#)
- [PERF01-BP04 Valuta l'impatto dei compromessi sui clienti e sull'efficienza dell'architettura](#)
- [PERF01-BP05 Usa politiche e architetture di riferimento](#)
- [PERF01-BP06 Uso del benchmarking per guidare le decisioni sull'architettura](#)
- [PERF01-BP07 Usa un approccio basato sui dati per le scelte architettoniche](#)

PERF01-BP01 Scopri e comprendi i servizi e le funzionalità cloud disponibili

Informati continuamente e identifica i servizi e le configurazioni disponibili che ti aiutano a prendere le decisioni giuste sull'architettura e a migliorare l'efficienza delle prestazioni dei carichi di lavoro.

Anti-pattern comuni:

- Utilizzi il cloud come data center in co-location.
- Non stai modernizzando la tua applicazione con la migrazione al cloud.

- Stai solo usando un tipo di archiviazione per tutte le cose che devono essere conservate in modo persistente.
- Se necessario, utilizzi tipi di istanze strettamente correlate ai tuoi standard attuali, ma più grandi.
- Distribuisci e gestisci le tecnologie disponibili come servizi gestiti.

Vantaggi dell'adozione di questa best practice: prendendo in considerazione nuovi servizi e configurazioni, puoi migliorare notevolmente le prestazioni, ridurre i costi e ottimizzare le attività necessarie per mantenere il carico di lavoro. Può anche aiutarti ad accelerare l'adozione di prodotti abilitati al cloud time-to-value.

Livello di rischio associato se questa best practice non fosse adottata: elevato

Guida all'implementazione

AWS rilascia continuamente nuovi servizi e funzionalità in grado di migliorare le prestazioni e ridurre il costo dei carichi di lavoro cloud. Utilizzare up-to-date questi nuovi servizi e funzionalità è fondamentale per mantenere l'efficacia delle prestazioni nel cloud. La modernizzazione dell'architettura dei carichi di lavoro consente inoltre di accelerare la produttività, promuovere l'innovazione e sbloccare ulteriori opportunità di crescita.

Passaggi dell'implementazione

- Esegui l'inventario del software e dell'architettura del carico di lavoro per i servizi correlati. Determina su quale categoria di prodotti ottenere ulteriori informazioni.
- Esplora AWS le offerte per identificare e conoscere i servizi e le opzioni di configurazione pertinenti che possono aiutarti a migliorare le prestazioni e ridurre i costi e la complessità operativa.
 - [Amazon Web Services Cloud](#)
 - [AWS Accademia](#)
 - [Cosa c'è di nuovo con AWS?](#)
 - [AWS Blog](#)
 - [AWS Skill Builder](#)
 - [AWS Eventi e webinar](#)
 - [AWS Training e certificazioni](#)
 - [AWS Canale Youtube](#)
 - [AWS Workshop](#)

- [Community AWS](#)
- Usa [Amazon Q](#) per ricevere informazioni e consigli pertinenti sui servizi.
- Usa gli ambienti sandbox non di produzione per comprendere e sperimentare nuovi servizi senza incorrere in costi aggiuntivi.
- Scopri servizi e funzionalità cloud sempre nuovi.

Risorse

Documenti correlati:

- [Overview of Amazon Web Services](#)
- [EC2Funzionalità di Amazon](#)
- [Impara step-by-step con un piano formativo per i AWS partner](#)
- [AWS Formazione e certificazione](#)
- [Il mio percorso di apprendimento per diventare un architetto di AWS soluzioni](#)
- [AWS Centro di architettura](#)
- [AWS Partner Network](#)
- [AWS Libreria di soluzioni](#)
- [AWS Centro di conoscenza](#)
- [Crea applicazioni moderne su AWS](#)

Video correlati:

- [AWS re:Invent 2023 - Cosa c'è di nuovo con Amazon EC2](#)
- [AWS re:Invent 2022 - Riduci i costi operativi e di infrastruttura con Amazon ECS](#)
- [AWS re:Invent 2023 - Costruisci con l'efficienza, l'agilità e l'innovazione del cloud con AWS](#)
- [AWS re:Invent 2022 - Implementa modelli ML per l'inferenza ad alte prestazioni e basso costo](#)
- [This is my Architecture](#)

Esempi correlati:

- [AWS Esempi](#)
- [AWS SDK Esempi](#)

PERF01-BP02 Utilizzo delle indicazioni del provider cloud o di un partner appropriato per conoscere gli schemi di architettura e le best practice

Usa le risorse aziendali del cloud come documentazione, solutions architect, servizi professionali o partner appropriati per guidare le tue decisioni sull'architettura. Queste risorse ti aiutano a rivedere e migliorare l'architettura per ottenere prestazioni ottimali.

Anti-pattern comuni:

- AWS è usato come un comune provider di servizi cloud.
- I servizi AWS vengono utilizzati in modo diverso rispetto alla loro progettazione iniziale.
- Le indicazioni vengono seguite senza considerare il contesto aziendale.

Vantaggi dell'adozione di questa best practice: avvalersi della guida di un provider di servizi cloud o di un partner appropriato può aiutarti a fare le scelte giuste per l'architettura del tuo carico di lavoro e darti fiducia nelle tue decisioni.

Livello di rischio associato se questa best practice non fosse adottata: medio

Guida all'implementazione

AWS offre un'ampia scelta di linee guida, documentazione e risorse che possono aiutarti a creare e gestire i carichi di lavoro del cloud in modo efficiente. La documentazione AWS fornisce esempi di codice, esercitazioni e spiegazioni dettagliate sui servizi. Oltre alla documentazione, AWS offre programmi di formazione e certificazione, solutions architect e servizi professionali che i clienti possono usare per esplorare diversi aspetti dei servizi cloud e implementare un'architettura cloud efficiente su AWS.

Sfrutta queste risorse per ottenere approfondimenti sulle informazioni e sulle best practice preziose per risparmiare tempo e ottenere risultati migliori nel Cloud AWS.

Passaggi dell'implementazione

- Consulta la documentazione e le linee guida AWS e segui le best practice. Queste risorse possono aiutarti a scegliere e configurare i servizi in modo efficace e a ottenere prestazioni migliori.
 - [Documentazione di AWS](#) (come guide utente e whitepaper)

- [Blog AWS](#)
- [AWS Training e certificazioni](#)
- [Canale YouTube di AWS](#)
- Partecipa agli eventi per i partner AWS (come summit AWS a livello mondiale, gruppi di utenti di AWS re:Invent e workshop) per apprendere dagli esperti AWS le best practice per l'utilizzo dei servizi AWS.
 - [Impara passo per passo con il Programma di apprendimento dei Partner AWS](#)
 - [Eventi e webinar AWS](#)
 - [Workshop AWS](#)
 - [Community AWS](#)
- Contatta AWS per ricevere assistenza quando ti occorrono ulteriori indicazioni o informazioni sui prodotti. AWS I Solutions Architect e i [servizi professionali di AWS](#) forniscono indicazioni per l'implementazione delle soluzioni. [AWS I partner](#) mettono a disposizione la propria conoscenza di AWS per aiutarti ad assicurare alla tua azienda agilità e innovazione.
- Usa [Supporto](#) se hai bisogno di supporto tecnico per utilizzare un servizio in modo efficace. I [nostri piani di supporto](#) sono pensati per offrirti il giusto mix di strumenti e competenze in modo da poter conseguire il successo con AWS ottimizzando le prestazioni, gestendo i rischi e tenendo sotto controllo i costi.

Risorse

Documenti correlati:

- [AWS Architecture Center](#)
- [AWS Partner Network](#)
- [Biblioteca di soluzioni di AWS](#)
- [Centro conoscenze di AWS](#)
- [Supporto AWS Enterprise](#)

Video correlati:

- [This is my Architecture](#)
- [AWS re:Invent 2023 - Advanced event-driven patterns with Amazon EventBridge](#)

- [AWS re:Invent 2023 - Implementing distributed design patterns on AWS](#)
- [AWS re:Invent 2023 - Application architecture as code](#)

Esempi correlati:

- [Esempi AWS](#)
- [Esempi di SDK AWS](#)
- [AWS Analytics Reference Architecture](#)

PERF01-BP03 Fattore di costo nelle decisioni architettoniche

Tieni conto dei costi nelle decisioni sull'architettura per migliorare l'utilizzo delle risorse e l'efficienza delle prestazioni del tuo carico di lavoro cloud. Quando si è consapevoli delle implicazioni dei costi del carico di lavoro cloud, è più probabile che si utilizzino risorse efficienti e si riducano le procedure inutili.

Anti-pattern comuni:

- Utilizzi una sola famiglia di istanze.
- Ometti di valutare le soluzioni con licenza rispetto alle soluzioni open-source.
- Non definisci le policy del ciclo di vita dell'archiviazione.
- Non recensisci i nuovi servizi e funzionalità di Cloud AWS
- Utilizzi solo lo storage a blocchi.

Vantaggi dell'adozione di questa best practice: la contabilizzazione dei costi nel processo decisionale consente di utilizzare risorse più efficienti ed esplorare altri investimenti.

Livello di rischio associato se questa best practice non fosse adottata: medio

Guida all'implementazione

L'ottimizzazione dei carichi di lavoro in base ai costi può migliorare l'utilizzo delle risorse ed evitare sprechi nel carico di lavoro cloud. Tenere conto dei costi nelle decisioni sull'architettura di solito include il corretto dimensionamento dei componenti del carico di lavoro e l'abilitazione dell'elasticità, comportando una migliore efficienza delle prestazioni del carico di lavoro cloud.

Passaggi dell'implementazione

- Stabilisci gli obiettivi di costo, come i limiti del budget, per il tuo carico di lavoro cloud.
- Identifica i componenti chiave, come istanze e archiviazione, che determinano il costo del carico di lavoro. Puoi usare [Calcolatore dei prezzi AWS](#) e [AWS Cost Explorer](#) per identificare i principali fattori di costo del carico di lavoro.
- Esamina i [modelli di prezzo](#) nel cloud, ad esempio istanze on-demand, riservate, Savings Plans e istanze spot.
- Segui le [best practice per l'ottimizzazione dei costi di Well-Architected](#) per ottimizzare questi componenti principali in termini di costi.
- Monitora e analizza continuamente i costi per identificare le opportunità di ottimizzazione dei costi nel tuo carico di lavoro.
 - Usa [Budget AWS](#) per ricevere gli avvisi per i costi inaccettabili.
 - Usa [AWS Compute Optimizer](#) o [AWS Trusted Advisor](#) per ottenere suggerimenti sull'ottimizzazione dei costi.
 - Usa [AWS Cost Anomaly Detection](#) per rilevare in modo automatico le anomalie dei costi e analizzare la causa principale.

Risorse

Documenti correlati:

- [Che cos'è AWS Billing and Cost Management?](#)
- [Ottimizzazione dei costi con AWS](#)
- [Scelta di una strategia di gestione dei AWS costi](#)
- [Una guida per principianti alla gestione AWS dei costi](#)
- [A Detailed Overview of the Cost Intelligence Dashboard](#)
- [AWS Architecture Center](#)
- [Biblioteca di soluzioni di AWS](#)
- [Centro conoscenze di AWS](#)

Video correlati:

- [This is my Architecture](#)

- [AWS re:Invent 2023 - Cosa c'è di nuovo con l'ottimizzazione dei costi AWS](#)
- [AWS re:Invent 2023 - Ottimizza costi e prestazioni e monitora i progressi verso la mitigazione](#)
- [AWS re:Invent 2023 - best practice per l'ottimizzazione dei costi di storage AWS](#)
- [AWS re:Invent 2023 - Ottimizza i costi nei tuoi ambienti con più account](#)

Esempi correlati:

- [AWS Compute Optimizer Codice demo](#)
- [Cost Optimization Workshop](#)
- [Cloud Financial Management Technical Implementation Playbooks](#)
- [Startup optimization: Tuning application performance for maximum efficiency](#)
- [Serverless Optimization Workshop \(Performance and Cost\)](#)
- [Scaling cost effective architectures](#)

PERF01-BP04 Valuta l'impatto dei compromessi sui clienti e sull'efficienza dell'architettura

Quando valuti i miglioramenti correlati alle prestazioni, determina quali scelte hanno impatto sui clienti e sull'efficienza del carico di lavoro. Ad esempio, se l'utilizzo di un datastore chiave-valore aumenta le prestazioni del sistema, è importante valutare in che modo la consistenza finale intrinseca di questo cambiamento avrà un impatto sui clienti.

Anti-pattern comuni:

- Ritieni che tutti i vantaggi prestazionali debbano essere implementati, anche se ci sono compromessi per l'implementazione.
- Valuti di apportare modifiche ai carichi di lavoro solo quando un problema prestazionale ha raggiunto un punto critico.

Vantaggi dell'adozione di questa best practice: quando si valutano potenziali miglioramenti relativi alle prestazioni, è necessario decidere se i compromessi per le modifiche sono accettabili con i requisiti del carico di lavoro. In alcuni casi, potrebbe essere necessario implementare controlli aggiuntivi per compensare i compromessi.

Livello di rischio associato se questa best practice non fosse adottata: elevato

Guida all'implementazione

Identifica le aree critiche della tua architettura in termini di prestazioni e impatto sui clienti. Stabilisci in che modo puoi apportare miglioramenti e quali compromessi comportano, oltre al loro impatto sul sistema e sull'esperienza degli utenti. L'implementazione di cache di dati, ad esempio, può contribuire a migliorare notevolmente le prestazioni ma richiede una strategia ben definita sulle modalità e sui tempi di aggiornamento o di invalidamento dei dati che vi sono contenuti, per evitare che il sistema si comporti in modo non corretto.

Passaggi dell'implementazione

- Comprendi i requisiti del tuo carico di SLAs lavoro e.
- Definisci chiaramente i fattori di valutazione. I fattori possono riguardare il costo, l'affidabilità, la sicurezza e le prestazioni del carico di lavoro.
- Seleziona l'architettura e i servizi in grado di soddisfare le tue esigenze.
- Conduci sperimentazioni e prove di fattibilità (POCs) per valutare i fattori di compromesso e l'impatto sui clienti e sull'efficienza dell'architettura. Di solito, i carichi di lavoro altamente disponibili, performanti e sicuri consumano più risorse cloud offrendo al contempo una esperienza cliente migliore. Comprendi i compromessi in termini di complessità, prestazioni e costi del tuo carico di lavoro. In genere, dare la priorità a due fattori va a scapito del terzo.

Risorse

Documenti correlati:

- [Amazon Builders' Library](#)
- [Amazon QuickSight KPIs](#)
- [Amazon CloudWatch RUM](#)
- [Documentazione di X-Ray](#)
- [Understand resiliency patterns and trade-offs to architect efficiently in the cloud](#)

Video correlati:

- [Ottimizza le applicazioni tramite Amazon CloudWatch RUM](#)
- [AWS re:Invent 2023 - Capacità, disponibilità, efficienza dei costi: scegline tre](#)

- [AWS re:Invent 2023 - Modelli di integrazione avanzati e compromessi per sistemi liberamente accoppiati](#)

Esempi correlati:

- [Misura il tempo di caricamento della pagina con Amazon CloudWatch Synthetics](#)
- [Client CloudWatch RUM Web Amazon](#)

PERF01-BP05 Usa politiche e architetture di riferimento

Utilizza le policy interne e le architetture di riferimento esistenti per la selezione dei servizi e delle configurazioni per una maggiore efficienza nella progettazione e nell'implementazione del carico di lavoro.

Anti-pattern comuni:

- Usi una vasta gamma di tecnologie che possono influire sul sovraccarico di gestione della tua azienda.

Vantaggi dell'adozione di questa best practice: la definizione di una policy per la scelta dell'architettura, della tecnologia e del fornitore consente di prendere decisioni rapidamente.

Livello di rischio associato se questa best practice non fosse adottata: medio

Guida all'implementazione

Avere policy interne nella selezione delle risorse e dell'architettura fornisce standard e linee guida da seguire quando si effettuano scelte architettoniche. Queste linee guida semplificano il processo decisionale nella scelta del servizio cloud giusto e possono contribuire a migliorare l'efficienza delle prestazioni. Implementi il carico di lavoro utilizzando policy o architetture di riferimento. Integra i servizi nell'implementazione cloud, quindi utilizza i test delle prestazioni per verificare che i requisiti prestazionali siano sempre rispettati.

Passaggi dell'implementazione

- Comprendi chiaramente i requisiti del tuo carico di lavoro cloud.
- Rivedi le policy interne ed esterne per identificare quelle più pertinenti.

- Utilizza le architetture di riferimento appropriate fornite dalle best practice AWS o di settore.
- Crea un contesto composto da policy, standard, architetture di riferimento e linee guida prescrittive per situazioni comuni. In questo modo i tuoi team possono muoversi più velocemente. Personalizza le risorse per il tuo settore verticale, se applicabile.
- Convalida queste policy e architetture di riferimento per il tuo carico di lavoro in ambienti sandbox.
- Resta up-to-date conforme agli standard e agli AWS aggiornamenti del settore per assicurarti che le tue policy e le architetture di riferimento contribuiscano a ottimizzare il carico di lavoro sul cloud.

Risorse

Documenti correlati:

- [AWS Architecture Center](#)
- [AWS Partner Network](#)
- [Biblioteca di soluzioni di AWS](#)
- [Centro conoscenze di AWS](#)
- [AWS Blog di architettura](#)

Video correlati:

- [This is my Architecture](#)
- [AWS re:Invent 2022 - Accelera il valore della tua azienda con SAP un'architettura di riferimento AWS](#)

Esempi correlati:

- [Esempi AWS](#)
- [AWS SDK Esempi](#)

PERF01-BP06 Uso del benchmarking per guidare le decisioni sull'architettura

Esegui il benchmark delle prestazioni di un carico di lavoro esistente per comprendere le prestazioni sul cloud e guidare le decisioni sull'architettura basate sui dati.

Anti-pattern comuni:

- Fai affidamento su valori di riferimento comuni che non sono indicativi delle caratteristiche del carico di lavoro.
- L'unico punto di riferimento è dato dal feedback e dalle percezioni dei clienti.

Vantaggi dell'adozione di questa best practice: misurazione dei miglioramenti in termini di prestazioni grazie al benchmarking dell'implementazione attuale.

Livello di rischio associato se questa best practice non fosse adottata: medio

Guida all'implementazione

Utilizza test sintetici di benchmarking per valutare le prestazioni dei componenti durante il carico di lavoro. Di solito, i benchmark sono più rapidi da configurare rispetto ai test di carico e vengono utilizzati per valutare la tecnologia di un componente specifico. Il benchmarking viene spesso utilizzato all'inizio di un nuovo progetto, quando non è ancora disponibile una soluzione completa da sottoporre a test di carico.

Puoi creare i tuoi test di benchmarking personalizzati oppure utilizzare test standard del settore, [come TPC-DS](#), per il benchmark dei carichi di lavoro. I benchmark di settore sono utili quando devi confrontare ambienti diversi. Quelli personalizzati, invece, sono indicati per analizzare tipi specifici di operazioni che prevedi di eseguire nell'architettura.

In fase di benchmarking, è importante effettuare delle operazioni preliminari sull'ambiente di test al fine di garantire la validità dei risultati. Dovrai eseguire lo stesso benchmark più volte, per verificare di avere acquisito ogni eventuale variazione nel corso del tempo.

Dal momento che, di solito, l'esecuzione dei benchmark è più rapida di quella dei test di carico, il benchmarking può essere utilizzato sin dalle prime fasi della pipeline di implementazione, così da fornire al team feedback più rapidi sulle deviazioni delle prestazioni. Quando valuti un cambiamento significativo in un componente o servizio, i benchmark possono essere un modo rapido per verificare se l'impegno necessario per apportare la modifica sia giustificato. L'utilizzo del benchmarking in combinazione con i test di carico è importante perché questi ultimi forniscono indicazioni sulle prestazioni del carico di lavoro in fase di produzione.

Passaggi dell'implementazione

- Pianifica e definisci:

- Definisci gli obiettivi, la baseline, gli scenari di test, le metriche, ad esempio l'utilizzo della CPU, la latenza o il throughput, e i KPI per il tuo benchmark.
- Concentrati sui requisiti degli utenti in termini di esperienza utente e su fattori come i tempi di risposta e l'accessibilità.
- Individua uno strumento di benchmark adatto al tuo carico di lavoro. Puoi utilizzare i servizi AWS (come [Amazon CloudWatch](#)) o uno strumento di terze parti compatibile con il tuo carico di lavoro.
- Configura ed esegui l'instrumentazione:
 - Imposta il tuo ambiente e configura le risorse.
 - Implementa il monitoraggio e la creazione di log per acquisire i risultati dei test.
- Esegui i test di benchmark e monitora:
 - Esegui i test di benchmark e monitora i parametri durante il test.
- Analizza e documenta:
 - Documenta il processo di benchmark e gli esiti.
 - Analizza i risultati per identificare i colli di bottiglia, le tendenze e le aree di miglioramento.
 - Usa i risultati dei test per prendere decisioni sull'architettura e modificare il carico di lavoro. Questa operazione può includere la modifica dei servizi o l'adozione di nuove funzionalità.
- Ottimizza e ripeti:
 - Modifica le configurazioni e le allocazioni delle risorse in base ai tuoi benchmark.
 - Ripeti il test del carico di lavoro dopo i cambiamenti per convalidare i miglioramenti.
 - Documenta le informazioni e ripeti il processo per identificare altre aree di miglioramento.

Risorse

Documenti correlati:

- [AWS Architecture Center](#)
- [AWS Partner Network](#)
- [Biblioteca di soluzioni di AWS](#)
- [Centro conoscenze di AWS](#)
- [Amazon CloudWatch RUM](#)
- [Amazon CloudWatch Synthetics](#)
- [Genomics workflows, Part 5: automated benchmarking](#)

- [Benchmark and optimize endpoint deployment in Amazon SageMaker JumpStart](#)

Video correlati:

- [AWS re:Invent 2023 - Benchmarking AWS Lambda cold starts](#)
- [Benchmarking stateful services in the cloud](#)
- [This is my Architecture](#)
- [Optimize applications through Amazon CloudWatch RUM](#)
- [Demo of Amazon CloudWatch Synthetics](#)

Esempi correlati:

- [Esempi AWS](#)
- [Esempi di SDK AWS](#)
- [Test del carico distribuito](#)
- [Misurazione dei tempi di caricamento delle pagine con Amazon CloudWatch Synthetics](#)
- [Client Web Amazon CloudWatch RUM](#)

PERF01-BP07 Usa un approccio basato sui dati per le scelte architettoniche

Definisci un approccio chiaro e basato sui dati per le scelte dell'architettura e verificare che vengano utilizzati i servizi e le configurazioni cloud corretti per soddisfare le tue esigenze aziendali specifiche.

Anti-pattern comuni:

- Ritieni che l'architettura corrente diventi statica e non venga aggiornata nel corso del tempo.
- Le tue scelte dell'architettura si basano su ipotesi e supposizioni.
- Introduci modifiche all'architettura nel tempo senza giustificazioni.

Vantaggi dell'adozione di questa best practice: con un approccio ben definito per le scelte dell'architettura, utilizzi i dati per influenzare la progettazione del carico di lavoro e prendere decisioni informate nel tempo.

Livello di rischio associato se questa best practice non fosse adottata: medio

Guida all'implementazione

Affidati all'esperienza e alle competenze interne in materia di cloud o utilizza risorse esterne, come casi d'uso pubblicati o whitepaper, per scegliere risorse e servizi per la tua architettura. È necessario definire con cura un processo che incoraggi la sperimentazione e il benchmarking con i servizi che possono essere utilizzati nel carico di lavoro.

I backlog dei carichi di lavoro critici devono consistere non solo in storie che offrono funzionalità rilevanti per l'azienda e gli utenti, ma anche in storie tecniche che definiscono la presentazione dell'architettura per il carico di lavoro. Questa presentazione include i nuovi progressi tecnologici e i nuovi servizi e li adotta sulla base di dati e giustificazioni adeguate. Verifica che l'architettura sia a prova di futuro e non diventi obsoleta.

Passaggi dell'implementazione

- Interagisci con le principali parti interessate per definire i requisiti del carico di lavoro, comprese le prestazioni, la disponibilità e le considerazioni sui costi. Includi fattori quali il numero di utenti e il modello di utilizzo del tuo carico di lavoro.
- Crea una presentazione dell'architettura o un backlog tecnologico a cui venga assegnata la priorità insieme al backlog funzionale.
- Valuta e identifica i diversi servizi cloud (per ulteriori dettagli, consulta [PERF01-BP01 Scopri e comprendi i servizi e le funzionalità cloud disponibili](#)).
- Esplora i diversi modelli di architettura, come microservizi o serverless, che soddisfano i tuoi requisiti di prestazioni (per maggiori dettagli, consulta [PERF01-BP02 Utilizzo delle indicazioni del provider cloud o di un partner appropriato per conoscere gli schemi di architettura e le best practice](#)).
- Consulta altri team, diagrammi di architettura e risorse, come AWS Solution Architects, [AWS Architecture Center](#) e [AWS Partner Network](#), per aiutarti a scegliere l'architettura giusta per il tuo carico di lavoro.
- Definisci i parametri, come il throughput e il tempo di risposta, che possono aiutarti a valutare le prestazioni del tuo carico di lavoro.
- Sperimenta e utilizza i parametri definiti per convalidare le prestazioni dell'architettura selezionata.

- Monitora continuamente e apporta le modifiche necessarie per mantenere ottimali le prestazioni della tua architettura.
- Documenta l'architettura e le decisioni selezionate come riferimento per aggiornamenti e apprendimenti futuri.
- Rivedi e aggiorna continuamente l'approccio di selezione dell'architettura in base agli apprendimenti, alle nuove tecnologie e ai parametri che indicano un problema o un cambiamento necessario nell'approccio attuale.

Risorse

Documenti correlati:

- [Biblioteca di soluzioni di AWS](#)
- [Centro conoscenze di AWS](#)
- [Modelli architettonici su cui creare applicazioni basate End-to-End sui dati AWS](#)

Video correlati:

- [This is my Architecture](#)
- [AWS re:Invent 2021 - Impresa basata sui dati: passare dalla visione al valore](#)
- [AWS re:Invent 2022 - Fornire architetture sostenibili e ad alte prestazioni](#)
- [AWS re:Invent 2023 - Ottimizza costi e prestazioni e monitora i progressi verso la mitigazione](#)
- [AWS re:Invent 2022 - AWS ottimizzazione: misure attuabili per risultati immediati](#)

Esempi correlati:

- [Esempi AWS](#)
- [AWS SDK Esempi](#)

Calcolo e hardware

La soluzione ottimale in termini di calcolo per un determinato carico di lavoro potrebbe variare in base alla progettazione dell'applicazione, ai modelli di utilizzo e alle impostazioni di configurazione. Le architetture possono utilizzare diverse soluzioni di calcolo per vari componenti e impiegare funzionalità diverse per migliorare le prestazioni. Selezionare la soluzione di calcolo sbagliata per un'architettura può ridurre l'efficienza delle prestazioni.

Questa area di interesse offre linee guida e best practice su come identificare e ottimizzare le opzioni di calcolo al fine di ottenere prestazioni di calcolo nel cloud efficienti.

Best practice

- [PERF02-BP01 Seleziona le migliori opzioni di elaborazione per il tuo carico di lavoro](#)
- [PERF02-BP02 Comprendi la configurazione e le funzionalità di elaborazione disponibili](#)
- [PERF02-BP03 Raccogli metriche relative al calcolo](#)
- [PERF02-BP04 Configurazione e dimensionamento corretto delle risorse di elaborazione](#)
- [PERF02-BP05 Scala le tue risorse di elaborazione in modo dinamico](#)
- [PERF02-BP06 Uso di acceleratori di elaborazione ottimizzati basati su hardware](#)

PERF02-BP01 Seleziona le migliori opzioni di elaborazione per il tuo carico di lavoro

La selezione dell'opzione di elaborazione più appropriata per il carico di lavoro consente di migliorare le prestazioni, ridurre i costi non necessari dell'infrastruttura e diminuire le attività operative richieste per mantenere il carico di lavoro.

Anti-pattern comuni:

- Si utilizza la stessa opzione di elaborazione utilizzata on-premises.
- Non si conoscono le opzioni, le funzionalità e le soluzioni di cloud computing e come queste migliorino le prestazioni di elaborazione.
- Si effettua il provisioning eccessivo dell'opzione di elaborazione per soddisfare i requisiti di dimensionamento o prestazioni, quando il passaggio a una nuova opzione di elaborazione soddisferebbe le caratteristiche del carico di lavoro in modo più preciso.

Vantaggi dell'adozione di questa best practice: identificando i requisiti di elaborazione e valutando le opzioni disponibili è possibile rendere il carico di lavoro più efficiente in termini di risorse.

Livello di rischio associato se questa best practice non fosse adottata: elevato

Guida all'implementazione

Per ottimizzare i carichi di lavoro cloud per l'efficienza delle prestazioni, è importante selezionare le opzioni di elaborazione più appropriate per il caso d'uso e i requisiti di prestazioni. AWS offre una varietà di opzioni di elaborazione che soddisfano diversi carichi di lavoro nel cloud. Ad esempio, puoi utilizzare [Amazon EC2](#) per avviare e gestire server virtuali, [AWS Lambda](#) eseguire codice senza dover fornire o gestire server, [Amazon ECS](#) o [Amazon EKS](#) per eseguire e gestire contenitori o [AWS Batch](#) elaborare grandi volumi di dati in parallelo. In base alle tue esigenze di dimensionamento ed elaborazione, scegli e configura la soluzione di elaborazione ottimale per la tua situazione. Puoi anche prendere in considerazione l'utilizzo di più tipi di soluzioni di elaborazione in un unico carico di lavoro in quanto ognuna ha i suoi vantaggi e svantaggi.

I passaggi seguenti ti guidano nella selezione delle opzioni di elaborazione giuste per soddisfare le caratteristiche del carico di lavoro e i requisiti prestazionali.

Passaggi dell'implementazione

- Comprendi i requisiti di elaborazione del tuo carico di lavoro. I requisiti essenziali da considerare includono le esigenze di elaborazione, gli schemi di traffico, gli schemi di accesso ai dati, le esigenze di dimensionamento e i requisiti di latenza.
- Scopri i vari [servizi di elaborazione AWS](#) per il tuo carico di lavoro. Per ulteriori informazioni, consulta [PERF01-BP01 Scopri e comprendi i servizi e le funzionalità cloud disponibili](#). Ecco alcune importanti opzioni di elaborazione AWS , le caratteristiche e i casi d'uso più comuni:

AWS servizio	Caratteristiche chiave	Casi di utilizzo comune
Amazon Elastic Compute Cloud (AmazonEC2)	Dispone di un'opzione dedicata per hardware, requisiti di licenza, ampia selezione di diverse famiglie di istanze, tipi di processori e acceleratori di elaborazione	Migrazioni con rehosting (lift and shift), applicazione monolitica, ambienti ibridi, applicazioni aziendali

AWS servizio	Caratteristiche chiave	Casi di utilizzo comune
Amazon Elastic Container Service (AmazonECS) , Amazon Elastic Kubernetes Service (Amazon) EKS	Implementazione semplice, ambienti coerenti, scalabile	Microservizi, ambienti ibridi
AWS Lambda	Servizio di elaborazione serverless che esegue il codice in risposta agli eventi e gestisce automaticamente le risorse di elaborazione sottostanti.	Microservizi, applicazioni basate su eventi
AWS Batch	Esegue il provisioning e la scalabilità in modo efficiente e dinamico di Amazon Elastic Container Service (AmazonECS) , Amazon Elastic Kubernetes Service EKS (Amazon AWS Fargate) e risorse di calcolo, con la possibilità di utilizzare istanze On-Demand o Spot in base alle tue esigenze lavorative	HPC, addestra modelli ML
Amazon Lightsail	Applicazione Linux e Windows preconfigurata per l'esecuzione di piccoli carichi di lavoro	Applicazioni Web semplici, sito Web personalizzato

- Valuta i costi (come la tariffa oraria o il trasferimento dei dati) e il sovraccarico di gestione (come l'applicazione di patch e il dimensionamento) associati a ciascuna opzione di elaborazione.
- Esegui esperimenti e benchmarking in un ambiente non di produzione per identificare quale opzione di elaborazione può soddisfare al meglio i requisiti del tuo carico di lavoro.
- Dopo aver sperimentato e identificato la tua nuova soluzione di calcolo, pianifica la migrazione e convalida i parametri prestazionali.

- Utilizza strumenti di AWS monitoraggio come [Amazon CloudWatch](#) e servizi di ottimizzazione [AWS Compute Optimizer](#) per ottimizzare continuamente le risorse di elaborazione in base a modelli di utilizzo reali.

Risorse

Documenti correlati:

- [Elaborazione in cloud con AWS](#)
- [Tipi di EC2 istanze Amazon](#)
- [Amazon EKS Containers: Amazon EKS Worker Nodes](#)
- [Amazon ECS Containers: istanze di Amazon ECS Container](#)
- [Funzioni: configurazione della funzione Lambda](#)
- [Prescriptive Guidance for Containers](#)
- [Prescriptive Guidance for Serverless](#)

Video correlati:

- [AWS re:Invent 2023 - AWS Graviton: il miglior rapporto prezzo/prestazioni per i tuoi carichi di lavoro AWS](#)
- [AWS re:Invent 2023 - Nuove funzionalità di intelligenza artificiale generativa di Amazon Elastic Compute Cloud in AMS](#)
- [AWS re:Invent 2023 - What's new with Amazon Elastic Compute Cloud](#)
- [AWS re:Invent 2023 - Smart savings: Amazon Elastic Compute Cloud cost-optimization strategies](#)
- [AWS re:Invent 2021 - Powering next-gen Amazon Elastic Compute Cloud: Deep dive on the Nitro System](#)
- [AWS re:Invent 2019 - Ottimizza prestazioni e costi per le tue risorse di calcolo AWS](#)
- [AWS re:Invent 2019 - Amazon Elastic Compute Cloud foundations](#)
- [AWS re:Invent 2022 - Implementa modelli ML per l'inferenza ad alte prestazioni e a basso costo](#)
- [AWS re:Invent 2019 - Ottimizza le prestazioni e i costi per il tuo calcolo AWS](#)
- [EC2Fondazioni Amazon](#)
- [Deploy ML models for inference at high performance and low cost](#)

Esempi correlati:

- [Migrating the Web application to containers](#)
- [Esecuzione di un "Hello, World!" serverless](#)
- [EKSWorkshop Amazon](#)
- [EC2Workshop Amazon](#)
- [Efficient and Resilient Workloads with Amazon Elastic Compute Cloud Auto Scaling](#)
- [Migrazione a AWS Graviton con Container Services](#)

PERF02-BP02 Comprendi la configurazione e le funzionalità di elaborazione disponibili

Comprendi le opzioni e le funzionalità di configurazione disponibili per il tuo servizio di calcolo in modo da fornire la giusta quantità di risorse e migliorare l'efficienza delle prestazioni.

Anti-pattern comuni:

- Non valuti le opzioni di calcolo o le famiglie di istanze disponibili rispetto alle caratteristiche del carico di lavoro.
- Esegui il provisioning eccessivo delle risorse di calcolo per soddisfare i requisiti di picco della domanda.

Vantaggi derivanti dall'adozione di questa best practice: acquisisci familiarità con le funzionalità e le configurazioni di AWS elaborazione in modo da poter utilizzare una soluzione di elaborazione ottimizzata per soddisfare le caratteristiche e le esigenze del carico di lavoro.

Livello di rischio associato se questa best practice non fosse adottata: medio

Guida all'implementazione

Ogni soluzione di calcolo ha disponibili configurazioni e funzionalità specifiche per supportare caratteristiche e requisiti diversi del carico di lavoro. Scopri in che modo puoi completare al meglio il tuo carico di lavoro e quali opzioni di configurazione sono le migliori per la tua applicazione. Esempi di queste opzioni includono famiglia di istanze, dimensioni, caratteristiche (, I/O)GPU, bursting, timeout, dimensioni delle funzioni, istanze di container e concorrenza. Se il tuo carico di lavoro utilizza la stessa opzione di elaborazione da più di quattro settimane e prevedi che le caratteristiche

rimarranno le stesse in futuro, puoi usarla [AWS Compute Optimizer](#) per scoprire se l'opzione di elaborazione attuale è adatta ai carichi di lavoro e dal punto di vista della memoria. CPU

Passaggi dell'implementazione

- Comprendi i requisiti del carico di lavoro (come CPU necessità, memoria e latenza).
- AWS Consulta la documentazione e le best practice per conoscere le opzioni di configurazione consigliate che possono contribuire a migliorare le prestazioni di elaborazione. Ecco alcune opzioni di configurazione chiave da considerare:

Opzione di configurazione	Esempi
Tipo di istanza	<ul style="list-style-type: none"> • Le istanze ottimizzate per il calcolo sono ideali per i carichi di lavoro che richiedono un rapporto v/memoria elevato e più elevato. CPU • Le istanze ottimizzate per la memoria offrono grandi quantità di memoria per carichi di lavoro intensivi in questo senso. • Le istanze ottimizzate per lo storage sono progettate per carichi di lavoro che richiedono un accesso sequenziale elevato in lettura e scrittura () allo storage locale. IOPS
Modello tariffario	<ul style="list-style-type: none"> • Le istanza on demand ti consentono di utilizzare la capacità di calcolo su base oraria o al secondo, senza impegni a lungo termine e sono ideali per il bursting oltre le esigenze di base per le prestazioni. • Savings Plans offrono risparmi significativi rispetto alle istanze on demand in cambio dell'impegno a utilizzare una quantità specifica di potenza di elaborazione per un periodo di uno o tre anni.

Opzione di configurazione	Esempi
	<ul style="list-style-type: none">• Le istanze spot ti consentono di sfruttare la capacità inutilizzata delle istanze con uno sconto per i carichi di lavoro stateless e tolleranti ai guasti.
Auto Scaling	Usa la configurazione Auto Scaling per abbinare le risorse di calcolo ai modelli di traffico.
Dimensionamento	<ul style="list-style-type: none">• Usa Compute Optimizer per ricevere un efficace suggerimento di machine learning riguardo alla configurazione più adatta alle tue caratteristiche di elaborazione.• Usa AWS Lambda Power Tuning per selezionare la configurazione migliore per la tua funzione Lambda.
Acceleratori di calcolo basati su hardware	<ul style="list-style-type: none">• Le istanze di elaborazione accelerata eseguono funzioni come l'elaborazione grafica o la corrispondenza dei modelli di dati in modo più efficiente rispetto alle alternative basate su base. CPU• Per i carichi di lavoro di machine learning, sfrutta l'hardware appositamente progettato e specifico per il tuo carico di lavoro, come AWS Trainium, Inferentia e Amazon AWS EC2 DL1

Risorse

Documenti correlati:

- [Elaborazione in cloud con AWS](#)
- [Tipi di EC2 istanze Amazon](#)

- [Controllo dello stato del processore per la tua EC2 istanza Amazon](#)
- [Amazon EKS Containers: Amazon EKS Worker Nodes](#)
- [Amazon ECS Containers: istanze di Amazon ECS Container](#)
- [Funzioni: configurazione della funzione Lambda](#)

Video correlati:

- [AWS re:Invent 2023 — AWS Graviton: il miglior rapporto prezzo/prestazioni per i tuoi carichi di lavoro AWS](#)
- [AWS re:Invent 2023 — Nuove funzionalità di intelligenza artificiale EC2 generativa di Amazon in AWS Management Console](#)
- [AWS re:Invent 2023 — Cosa c'è di nuovo con Amazon EC2](#)
- [AWS re:Invent 2023 — Risparmio intelligente: strategie di ottimizzazione dei costi di Amazon EC2](#)
- [AWS re:Invent 2021 — Potenziamento della nuova generazione di EC2 Amazon: Deep dive on the Nitro System](#)
- [AWS re:Invent 2019 — Fondamenti Amazon EC2](#)
- [AWS re:Invent 2022 — Ottimizzazione di Amazon EKS per prestazioni e costi AWS](#)

Esempi correlati:

- [Codice dimostrativo di Compute Optimizer](#)
- [Workshop sulle istanze EC2 spot di Amazon](#)
- [Carichi di lavoro efficienti e resilienti con Amazon EC2 AWS Auto Scaling](#)
- [Workshop per sviluppatori Graviton](#)
- [AWS per la giornata di immersione dei carichi di lavoro Microsoft](#)
- [AWS per una giornata di immersione nei carichi di lavoro Linux](#)
- [AWS Compute Optimizer Codice dimostrativo](#)
- [EKSWorkshop Amazon](#)

PERF02-BP03 Raccogli metriche relative al calcolo

Registra e monitora i parametri relativi all'elaborazione per comprendere meglio le prestazioni delle tue risorse di elaborazione e migliorarne le prestazioni e l'utilizzo.

Anti-pattern comuni:

- Utilizzi solo i file di log manuali per la ricerca dei parametri.
- Utilizzi solo i parametri predefiniti registrati dal software di monitoraggio.
- Revisione dei parametri solo quando c'è un problema.

Vantaggi dell'adozione di questa best practice: la raccolta dei parametri relativi alle prestazioni ti aiuta ad allineare le prestazioni delle applicazioni ai requisiti aziendali per garantire il rispetto delle esigenze dei carichi di lavoro. Può anche aiutarti a migliorare costantemente le prestazioni e l'utilizzo delle risorse del tuo carico di lavoro.

Livello di rischio associato se questa best practice non fosse adottata: elevato

Guida all'implementazione

I carichi di lavoro del cloud possono generare grandi volumi di dati quali parametri, log ed eventi. Nel Cloud AWS, la raccolta delle metriche è un passaggio fondamentale per migliorare la sicurezza, l'efficienza dei costi, le prestazioni e la sostenibilità. AWS fornisce un'ampia gamma di metriche relative alle prestazioni utilizzando servizi di monitoraggio come [Amazon CloudWatch](#) per fornirti informazioni preziose. Metriche come CPU l'utilizzo, l'utilizzo della memoria, l'I/O del disco e la rete in entrata e in uscita possono fornire informazioni sui livelli di utilizzo o sui colli di bottiglia delle prestazioni. Utilizza tali parametri come parte di un approccio basato sui dati per ottimizzare e ottimizzare le risorse del tuo carico di lavoro. L'ideale sarebbe raccogliere tutti i parametri relativi alle tue risorse di elaborazione in un'unica piattaforma con policy di conservazione implementate per supportare costi e obiettivi operativi.

Passaggi dell'implementazione

- Identifica quali parametri relativi alle prestazioni sono rilevanti per il tuo carico di lavoro. Raccogli i parametri sull'utilizzo delle risorse e sul modo in cui opera il tuo carico di lavoro nel cloud (come il tempo di risposta e il throughput).
 - [Metriche EC2 predefinite di Amazon](#)

- [Metriche ECS predefinite di Amazon](#)
- [Metriche EKS predefinite di Amazon](#)
- [Parametri predefiniti di Lambda](#)
- [Parametri EC2 della memoria e del disco di Amazon](#)
- Scegli e configura la soluzione di registrazione e monitoraggio giusta per il tuo carico di lavoro.
 - [AWS native Observability](#)
 - [AWS Distro per OpenTelemetry](#)
 - [Amazon Managed Service per Prometheus](#)
- Definisci il filtro e l'aggregazione richiesti per i parametri in base ai requisiti del tuo carico di lavoro.
 - [Quantifica i parametri delle applicazioni personalizzate con Amazon CloudWatch Logs e filtri metrici](#)
 - [Raccogli metriche personalizzate con il tagging CloudWatch strategico di Amazon](#)
- Configura le policy di conservazione dei dati per i parametri in modo che corrispondano ai tuoi obiettivi operativi e di sicurezza.
 - [Conservazione dei dati predefinita per le metriche CloudWatch](#)
 - [Conservazione dei dati predefinita per i registri CloudWatch](#)
- Se necessario, crea allarmi e notifiche per i parametri in modo da rispondere in modo proattivo ai problemi relativi alle prestazioni.
 - [Crea allarmi per metriche personalizzate utilizzando il rilevamento delle anomalie di Amazon CloudWatch](#)
 - [Crea metriche e allarmi per pagine Web specifiche con Amazon CloudWatch RUM](#)
- Usa l'automazione per implementare gli agenti di aggregazione di parametri e log.
 - [AWS Systems Manager automazione](#)
 - [OpenTelemetryCollezionista](#)

Risorse

Documenti correlati:

- [Monitoraggio e osservabilità](#)
- [Migliori pratiche: implementazione dell'osservabilità con AWS](#)
- [CloudWatch Documentazione Amazon](#)

- [Raccogli metriche e log EC2 dalle istanze Amazon e dai server locali con l'agente CloudWatch](#)
- [Accesso ad Amazon CloudWatch Logs per AWS Lambda](#)
- [Utilizzo dei CloudWatch log con istanze di container](#)
- [Publish custom metrics](#)
- [AWS Answers: Centralized Logging](#)
- [AWS Servizi che pubblicano metriche CloudWatch](#)
- [Monitoraggio di Amazon EKS su AWS Fargate](#)

Video correlati:

- [AWS re:Invent 2023 — \[LAUNCH\] Monitoraggio delle applicazioni per carichi di lavoro moderni](#)
- [AWS re:Invent 2023 — Implementazione dell'osservabilità delle applicazioni](#)
- [AWS re:Invent 2023 — Creazione di una strategia di osservabilità efficace](#)
- [AWS re:Invent 2023 — Osservabilità senza interruzioni con Distro per AWS OpenTelemetry](#)
- [Gestione delle prestazioni delle applicazioni su AWS](#)

Esempi correlati:

- [AWS per Linux Workload Immersion Day- Amazon CloudWatch](#)
- [Monitoraggio di ECS cluster e container Amazon](#)
- [Monitoraggio con CloudWatch dashboard Amazon](#)
- [EKSWorkshop Amazon](#)

PERF02-BP04 Configurazione e dimensionamento corretto delle risorse di elaborazione

Configura e dimensiona correttamente le risorse di elaborazione per soddisfare i requisiti di prestazioni del carico di lavoro ed evitare un utilizzo insufficiente o eccessivo delle risorse.

Anti-pattern comuni:

- Ignori i requisiti di prestazioni del carico di lavoro, con il risultato del provisioning eccessivo o insufficiente delle risorse di elaborazione.

- Scegli semplicemente l'istanza più grande o più piccola disponibile per tutti i carichi di lavoro.
- Usi una sola famiglia di istanze per semplificare la gestione.
- Ignori i consigli di Compute AWS Cost Explorer Optimizer o di Compute Optimizer per il corretto dimensionamento.
- Non rivaluti il carico di lavoro in base all'idoneità dei nuovi tipi di istanza.
- Certifici solo un numero limitato di configurazioni di istanza per l'organizzazione.

Vantaggi dell'adozione di questa best practice il corretto dimensionamento delle risorse di elaborazione garantisce un funzionamento ottimale nel cloud evitando il provisioning eccessivo o insufficiente delle risorse. Il corretto dimensionamento delle risorse di elaborazione comporta in genere prestazioni ottimali e una migliore esperienza cliente, riducendo al contempo i costi.

Livello di rischio associato se questa best practice non fosse adottata: medio

Guida all'implementazione

Il dimensionamento corretto consente alle organizzazioni di gestire la propria infrastruttura cloud in modo efficiente ed economico, rispettando al contempo le esigenze aziendali. L'eccessivo provisioning delle risorse cloud può comportare costi aggiuntivi, mentre un approvvigionamento insufficiente può comportare prestazioni scadenti e un'esperienza negativa per il cliente. AWS fornisce strumenti come [AWS Compute Optimizer](#) e [AWS Trusted Advisor](#) che utilizzano dati storici per fornire consigli sul corretto dimensionamento delle risorse di elaborazione.

Passaggi dell'implementazione

- Scegli il tipo di istanza più adatto alle tue esigenze:
 - [Come faccio a scegliere il tipo di EC2 istanza Amazon appropriato per il mio carico di lavoro?](#)
 - [Selezione del tipo di istanza basata sugli attributi per Amazon Fleet EC2](#)
 - [Create an Auto Scaling group using attribute-based instance type selection](#)
 - [Optimizing your Kubernetes compute costs with Karpenter consolidation](#)
- Analizza le varie caratteristiche prestazionali del tuo carico di lavoro e come queste caratteristiche si relazionano alla memoria, alla rete e all'utilizzo. CPU Utilizza questi dati per scegliere le risorse che meglio corrispondono al profilo del tuo carico di lavoro e agli obiettivi di prestazioni.
- Monitora l'utilizzo delle risorse utilizzando strumenti di AWS monitoraggio come Amazon CloudWatch.

- Seleziona la configurazione corretta per la risorsa di elaborazione.
 - Per carichi di lavoro temporanei, valuta i [CloudWatch parametri di Amazon](#) dell'istanza, ad esempio `CPUUtilization` per identificare se l'istanza è sottoutilizzata o sovrautilizzata.
 - Per carichi di lavoro stabili, controlla gli strumenti di AWS corretto dimensionamento, ad esempio e a intervalli regolari, per identificare le opportunità di ottimizzazione e dimensionamento corretto della risorsa di elaborazione AWS Compute Optimizer. AWS Trusted Advisor
- Esegui il test delle modifiche apportate alla configurazione in un ambiente non di produzione prima di implementarle in un ambiente live.
- Rivaluta costantemente nuove offerte di elaborazione e confrontale con le esigenze del carico di lavoro.

Risorse

Documenti correlati:

- [Cloud Compute con AWS](#)
- [Tipi di EC2 istanze Amazon](#)
- [Amazon ECS Containers: istanze di Amazon ECS Container](#)
- [Amazon EKS Containers: Amazon EKS Worker Nodes](#)
- [Funzioni: configurazione della funzione Lambda](#)
- [Controllo dello stato del processore per la tua EC2 istanza Amazon](#)

Video correlati:

- [EC2Fondazioni Amazon](#)
- [AWS re:Invent 2023 — AWS Graviton: il miglior rapporto prezzo/prestazioni per i tuoi carichi di lavoro AWS](#)
- [AWS re:Invent 2023 — Nuove funzionalità di intelligenza artificiale EC2 generativa di Amazon in AWS Management Console](#)
- [AWS re:Invent 2023 — Cosa c'è di nuovo con Amazon EC2](#)
- [AWS re:Invent 2023 — Risparmio intelligente: strategie di ottimizzazione dei costi di Amazon EC2](#)
- [AWS re:Invent 2021 — Potenziamento della nuova generazione di EC2 Amazon: Deep dive on the Nitro System](#)

- [AWS re:Invent 2019 — Fondamenti Amazon EC2](#)

Esempi correlati:

- [AWS Compute Optimizer Codice demo](#)
- [EKSOfficina Amazon](#)
- [Right-sizing recommendations](#)

PERF02-BP05 Scala le tue risorse di elaborazione in modo dinamico

Sfrutta l'elasticità del cloud per scalare dinamicamente le risorse di elaborazione per soddisfare le tue esigenze ed evitare un provisioning eccessivo o insufficiente per il tuo carico di lavoro.

Anti-pattern comuni:

- Risposta agli allarmi aumentando manualmente la capacità.
- Utilizzi le stesse linee guida per il dimensionamento (generalmente infrastruttura statica) di quelle on-premises.
- Dopo un evento di dimensionamento, lasci una capacità aumentata anziché ridurre il dimensionamento.

Vantaggi dell'adozione di questa best practice: la configurazione e il test dell'elasticità delle risorse di elaborazione possono aiutarti a risparmiare denaro, mantenere i benchmark delle prestazioni e migliorare l'affidabilità al variare del traffico.

Livello di rischio associato se questa best practice non fosse adottata: elevato

Guida all'implementazione

AWS offre la flessibilità necessaria per aumentare o ridurre le risorse in modo dinamico attraverso una varietà di meccanismi di scalabilità al fine di soddisfare le variazioni della domanda. In combinazione con i parametri relativi all'elaborazione, il dimensionamento dinamico consente ai carichi di lavoro di rispondere automaticamente alle modifiche e utilizzare il set ottimale di risorse di elaborazione per raggiungere l'obiettivo.

Puoi adottare varie strategie di approccio per associare l'offerta di risorse alla domanda.

- Approccio al tracciamento degli obiettivi: monitora il parametro di dimensionamento e aumenta o diminuisci automaticamente la capacità in base alle esigenze.
- Dimensionamento predittivo: procedi a ridurre orizzontalmente in previsione delle tendenze giornaliere e settimanali.
- Approccio basato sulla pianificazione: imposta il tuo programma di dimensionamento in base alle variazioni di carico prevedibili.
- Scalabilità del servizio: scegli i servizi (come quelli serverless) che si dimensionano automaticamente per progettazione.

Assicurati che le implementazioni dei carichi di lavoro siano in grado di gestire eventi che prevedono l'aumentare verticalmente e il ridurre verticalmente.

Passaggi dell'implementazione

- Istanze di elaborazione, container e funzioni forniscono tutti meccanismi di elasticità, in combinazione con il dimensionamento automatico o sotto forma di funzionalità del servizio. Ecco alcuni esempi di meccanismi di dimensionamento automatico:

Meccanismo di scalabilità automatica	Dove usarlo
Amazon EC2 Auto Scaling	Per assicurarti di disporre del numero corretto di EC2 istanze Amazon disponibili per gestire il carico di utenti per la tua applicazione.
Application Auto Scaling	Per scalare automaticamente le risorse per singoli AWS servizi oltre Amazon, EC2 come AWS Lambda funzioni o servizi Amazon Elastic Container Service (AmazonECS) .
Kubernetes Cluster Autoscaler/Karpenter	Dimensiona automaticamente i cluster Kubernetes.

- La scalabilità viene spesso discussa in relazione a servizi di elaborazione come Amazon EC2 Instances o funzioni. AWS Lambda Assicurati di considerare anche la configurazione di servizi non di calcolo come [AWS Glue](#) per soddisfare la domanda.
- Verifica che i parametri per il dimensionamento corrispondano alle caratteristiche del carico di lavoro da implementare. Se stai implementando un'applicazione di transcodifica video, è previsto

un CPU utilizzo del 100% e non dovrebbe essere il parametro principale. Utilizza la profondità della coda dei processi di transcodifica. Se necessario, puoi utilizzare una [metrica personalizzata](#) per la tua policy di dimensionamento. Per scegliere le metriche giuste, consulta le seguenti linee guida per AmazonEC2:

- La metrica deve essere una metrica di utilizzo valida e descrivere il livello di impiego di un'istanza.
- Il valore del parametro deve aumentare e diminuire in proporzione al numero di istanze nel gruppo con scalabilità automatica.
- Assicurati di utilizzare il [dimensionamento dinamico](#) anziché il [dimensionamento manuale](#) per il tuo gruppo Auto Scaling. È consigliabile utilizzare le [policy di dimensionamento del monitoraggio degli obiettivi](#) nel dimensionamento dinamico
- Verifica che le implementazioni dei carichi di lavoro siano in grado di gestire entrambi gli eventi di dimensionamento (aumento e riduzione). Ad esempio, puoi usare la [cronologia delle attività](#) per verificare le attività di ridimensionamento per un gruppo Auto Scaling.
- Analizza il tuo carico di lavoro per individuare modelli prevedibili e dimensionare le tue risorse in modo proattivo, anticipando variazioni nella domanda previste e pianificate. Con il dimensionamento predittivo puoi eliminare la necessità di offrire capacità in eccedenza. Per maggiori dettagli, consulta [Predictive Scaling with Amazon Auto EC2 Scaling](#).

Risorse

Documenti correlati:

- [Cloud Compute con AWS](#)
- [Tipi di EC2 istanze Amazon](#)
- [Amazon ECS Containers: istanze di Amazon ECS Container](#)
- [Amazon EKS Containers: Amazon EKS Worker Nodes](#)
- [Funzioni: configurazione della funzione Lambda](#)
- [Controllo dello stato del processore per la tua EC2 istanza Amazon](#)
- [Approfondimento su Amazon ECS Cluster Auto Scaling](#)
- [Introducing Karpenter – An Open-Source High-Performance Kubernetes Cluster Autoscaler](#)

Video correlati:

- [AWS re:Invent 2023 — AWS Graviton: il miglior rapporto prezzo/prestazioni per i tuoi carichi di lavoro AWS](#)
- [AWS re:Invent 2023 — Nuove funzionalità di intelligenza artificiale EC2 generativa di Amazon nella console di gestione AWS](#)
- [AWS re:Invent 2023 — Cosa c'è di nuovo con Amazon EC2](#)
- [AWS re:Invent 2023 — Risparmio intelligente: strategie di ottimizzazione dei costi di Amazon EC2](#)
- [AWS re:Invent 2021 — Potenziamento della nuova generazione di EC2 Amazon: Deep dive on the Nitro System](#)
- [AWS re:Invent 2019 — Fondamenti Amazon EC2](#)

Esempi correlati:

- [Esempi di EC2 gruppi Amazon Auto Scaling](#)
- [EKWorkshop Amazon](#)
- [Scala i tuoi EKS carichi di lavoro Amazon eseguendoli su IPv6](#)

PERF02-BP06 Uso di acceleratori di elaborazione ottimizzati basati su hardware

Usa gli acceleratori hardware per eseguire determinate funzioni in modo più efficiente rispetto alle alternative basate sulla CPU.

Anti-pattern comuni:

- Nel carico di lavoro non hai confrontato un'istanza per uso generico con un'istanza dedicata in grado di offrire prestazioni più elevate e costi inferiori.
- Usi gli acceleratori di calcolo basati su hardware per attività in cui sono più efficienti le alternative basate su CPU.
- Utilizzo delle GPU non monitorato.

Vantaggi dell'adozione di questa best practice: utilizzando gli acceleratori basati su hardware, come le unità di elaborazione grafica (GPU) e le serie di porte programmabili sul campo (FPGA), è possibile eseguire determinate funzioni di elaborazione in modo più efficiente.

Livello di rischio associato se questa best practice non fosse adottata: medio

Guida all'implementazione

Le istanze a calcolo accelerato forniscono l'accesso agli acceleratori di calcolo basati su hardware, come GPU e FPGA. Questi acceleratori hardware eseguono alcune funzioni, come l'elaborazione grafica o la rilevazione della corrispondenza dei modelli di dati, in modo più efficiente rispetto alle alternative basate su CPU. Molti carichi di lavoro accelerati, come il rendering grafico, la transcodifica e il machine learning, sono altamente variabili in termini di utilizzo di risorse. Esegui questo hardware solo per il tempo necessario e disattivalo con l'automazione quando non serve per migliorare l'efficienza complessiva delle prestazioni.

Passaggi dell'implementazione

- Identifica le [istanza a calcolo accelerato](#) in grado di soddisfare i tuoi requisiti.
- Per i carichi di lavoro di machine learning, sfrutta l'hardware specifico per il tuo carico di lavoro, come [AWS Trainium](#), [AWS Inferentia](#) e [Amazon EC2 DL1](#). AWS Le istanze Inferentia come le istanze Inf2 [offrono fino al 50% in più di prestazioni per watt rispetto alle istanze Amazon EC2 paragonabili](#).
- Raccogli i parametri di utilizzo delle istanze a calcolo accelerato. Ad esempio, puoi utilizzare l'agente CloudWatch per acquisire metriche quali `utilization_gpu` e `utilization_memory` per le tue GPU, come illustrato in [Collect NVIDIA GPU metrics with Amazon CloudWatch](#).
- Ottimizza il codice, il funzionamento della rete e le impostazioni degli acceleratori hardware per garantire il pieno utilizzo dell'hardware sottostante.
 - [Ottimizza le impostazioni GPU](#)
 - [Monitoraggio e ottimizzazione delle GPU nell'AMI per il deep learning](#)
 - [Optimizing I/O for GPU performance tuning of deep learning training in Amazon SageMaker](#)
- Utilizza le librerie e i driver per GPU più recenti e performanti.
- Utilizza l'automazione per rilasciare le istanze GPU non in uso.

Risorse

Documenti correlati:

- [Utilizzo di GPU su Amazon Elastic Container Service](#)
- [Istanze GPU](#)
- [Istanze con AWS Trainium](#)

- [Istanze con AWS Inferentia](#)
- [Let's Architect! Architecting with custom chips and accelerators](#)

- [Calcolo accelerato](#)
- [Amazon EC2 VT1 Instances](#)
- [Come faccio a scegliere il tipo di istanza Amazon EC2 appropriata per il mio carico di lavoro?](#)
- [Choose the best AI accelerator and model compilation for computer vision inference with Amazon SageMaker](#)

Video correlati:

- AWS re:Invent 2021 - [How to select Amazon Elastic Compute Cloud GPU instances for deep learning](#)
- [AWS re:Invent 2022 - \[NEW LAUNCH!\] Introducing AWS Inferentia2-based Amazon EC2 Inf2 instances](#)
- [AWS re:Invent 2022 - Accelerate deep learning and innovate faster with AWS Trainium](#)
- [AWS re:Invent 2022 - Deep learning on AWS with NVIDIA: From training to deployment](#)

Esempi correlati:

- [Amazon SageMaker and NVIDIA GPU Cloud \(NGC\)](#)
- [Use SageMaker with Trainium and Inferentia for optimized deep learning training and inferencing workloads](#)
- [Optimizing NLP models with Amazon Elastic Compute Cloud Inf1 instances in Amazon SageMaker](#)

Gestione dei dati

La soluzione ottimale per la gestione dei dati in un sistema specifico varia in base al tipo di dati (blocco, file o oggetto), agli schemi di accesso (casuali o sequenziali), al throughput necessario, alla frequenza di accesso (online, offline, archivio), alla frequenza di aggiornamento (WORM, dinamico) e ai vincoli di disponibilità e durata. I carichi di lavoro Well-Architected utilizzano archivi dati appositamente progettati che impiegano diverse funzionalità per migliorare le prestazioni.

Quest'area di interesse offre linee guida e best practice per ottimizzare l'archiviazione dei dati, i modelli di spostamento e accesso e l'efficienza delle prestazioni dell'archiviazione di dati.

Best practice

- [PERF03-BP01 Utilizza un archivio dati appositamente progettato che supporti al meglio i requisiti di accesso e archiviazione dei dati](#)
- [PERF03-BP02 Valuta le opzioni di configurazione disponibili per l'archivio dati](#)
- [PERF03-BP03 Raccolta e registrazione delle metriche delle prestazioni degli archivi dati](#)
- [PERF03-BP04 Implementazione di strategie per migliorare le prestazioni delle query nel datastore](#)
- [PERF03-BP05 Implementare modelli di accesso ai dati che utilizzano la memorizzazione nella cache](#)

PERF03-BP01 Utilizza un archivio dati appositamente progettato che supporti al meglio i requisiti di accesso e archiviazione dei dati

Comprendi le caratteristiche dei dati (come la condivisione, le dimensioni, la dimensione della cache, gli schemi di accesso, la latenza, il throughput e la persistenza dei dati) per selezionare i data store (archiviazione o database) dedicati per il tuo carico di lavoro.

Anti-pattern comuni:

- Continui a utilizzare un datastore per via dell'esperienza e delle competenze interne relative a quel particolare tipo di soluzione di database.
- Ritieni che tutti i carichi di lavoro abbiano requisiti di accesso e archiviazione di dati simili.
- Non hai implementato un catalogo di dati per eseguire l'inventario dei tuoi asset.

Vantaggi dell'adozione di questa best practice: la comprensione delle caratteristiche e dei requisiti dei dati ti consente di determinare la tecnologia di archiviazione più efficiente e performante appropriata per le tue esigenze del carico di lavoro.

Livello di rischio associato se questa best practice non fosse adottata: elevato

Guida all'implementazione

Quando selezionate e implementate l'archiviazione dei dati, assicuratevi che le caratteristiche di interrogazione, scalabilità e archiviazione supportino i requisiti relativi ai dati del carico di lavoro. AWS offre numerose tecnologie di archiviazione dei dati e di database, tra cui storage a blocchi, storage di oggetti, storage in streaming, file system, database relazionali, chiave-valore, documentali, in memoria, grafici, di serie temporali e di registro. Ogni soluzione di gestione dei dati offre soluzioni e configurazioni adatte a gestire i tuoi casi d'uso e modelli di dati. Comprendendo le caratteristiche e i requisiti dei dati, è possibile abbandonare la tecnologia di storage monolitica e adottare approcci restrittivi per concentrarsi sulla gestione appropriata dei dati. one-size-fits-all

Passaggi dell'implementazione

- Esegui un inventario dei vari tipi di dati esistenti nel tuo carico di lavoro.
- Comprendi e documenta le caratteristiche e i requisiti dei dati, tra cui:
 - Tipo di dati (non strutturati, semi-strutturati, relazionali)
 - Volume e crescita dei dati
 - Durabilità dei dati: persistenti, effimeri, transitori
 - ACIDrequisiti (atomicità, coerenza, isolamento, durabilità)
 - Schemi di accesso ai dati (con uso intensivo di lettura o scrittura)
 - Latenza
 - Prestazioni
 - IOPS(operazioni di ingresso/uscita al secondo)
 - Periodo di conservazione dei dati
- Scopri i diversi archivi di dati (servizi [di archiviazione](#) e [database](#)) disponibili per il tuo carico di lavoro e AWS che possono soddisfare le caratteristiche dei tuoi dati, come descritto in. [PERF01-BP01 Scopri e comprendi i servizi e le funzionalità cloud disponibili](#) Alcuni esempi di tecnologie di archiviazione AWS e delle loro caratteristiche chiave sono:

Tipo	AWS Servizi	Caratteristiche chiave
Archiviazione di oggetti	Amazon S3	Scalabilità illimitata, alta disponibilità e molteplici opzioni di accessibilità. L'accesso a oggetti e il relativo trasferimento da e verso Amazon S3 può utilizzare un servizio, come Transfer Acceleration o Punti di accesso , per supportare la posizione, le esigenze di sicurezza e i modelli di accesso.
Archiviazione	Amazon S3 Glacier	Progettato per l'archiviazione dei dati.
Archiviazione in streaming	Amazon Kinesis Streaming gestito da Amazon per Apache Kafka (Amazon MSK)	Acquisizione e archiviazione efficienti dei dati in streaming.
File system condiviso	Amazon Elastic File System (AmazonEFS)	File system montabile a cui è possibile accedere da più tipi di soluzioni di calcolo.

Tipo	AWS Servizi	Caratteristiche chiave
File system condiviso	Amazon FSx	Basato sulle più recenti soluzioni di AWS elaborate per supportare quattro file system di uso comune: Open NetApp ONTAPZFS, Windows File Server e Lustre. FSx <u>La latenza, la velocità effettiva e la velocità effettiva</u> di Amazon IOPS variano in base al file system e devono essere prese in considerazione quando si seleziona il file system giusto per le esigenze di carico di lavoro.
Storage a blocchi	Amazon Elastic Block Store (AmazonEBS)	Servizio di storage a blocchi scalabile e ad alte prestazioni progettato per Amazon Elastic Compute Cloud (Amazon). EC2 Amazon EBS include storage SSD supportato per carichi di lavoro transazionali e intensivi e HDD storage supportato per carichi di lavoro con throughput IOPS intensivo.

Tipo	AWS Servizi	Caratteristiche chiave
Database relazionale	Amazon Aurora , AmazonRDS, Amazon Redshift .	Progettato per supportare transazioni ACID (atomicità, coerenza, isolamento, durabilità) e mantenere l'integrità referenziale e una forte coerenza dei dati. Molte applicazioni tradizionali, la pianificazione delle risorse aziendali (ERP), la gestione delle relazioni con i clienti (CRM) e l'e-commerce utilizzano database relazionali per archiviare i propri dati.
Database chiave-valore	Amazon DynamoDB	Ottimizzato per schemi di accesso di uso comune, in genere per archiviare e recuperare grandi volumi di dati. Le app Web dal traffico elevato, i sistemi di e-commerce e le applicazioni di videogiochi sono casi d'uso tipici dei database chiave-valore.
Database di documenti	Amazon DocumentDB	Progettato per archiviare dati semistrutturati come documenti similiJSON. Questi database aiutano gli sviluppatori a creare e aggiornare rapidamente applicazioni quali gestione di contenuti, cataloghi e profili utente.

Tipo	AWS Servizi	Caratteristiche chiave
Database in memoria	Amazon ElastiCache , Amazon MemoryDB per Redis	Vengono utilizzati per applicazioni che richiedono accesso in tempo reale ai dati, bassissima latenza ed elevatissimo throughput. È possibile utilizzare database in memoria per la memorizzazione nella cache delle applicazioni, la gestione delle sessioni, la classifica dei giochi, l'archivio delle caratteristiche ML a bassa latenza, il sistema di messaggistica dei microservizi e un meccanismo di streaming a elevato throughput.
Database a grafo	Amazon Neptune	Utilizzato con le applicazioni che devono navigare ed eseguire query su milioni di relazioni tra set di dati a grafo altamente connessi, con una latenza misurata in millisecondi su larga scala. Molte aziende utilizzano database a grafo per il rilevamento di attività fraudolente, i social network e i motori di raccomandazione.

Tipo	AWS Servizi	Caratteristiche chiave
Database di serie temporali	Amazon Timestream	Utilizzato per raccogliere, sintetizzare e derivare in modo efficiente approfondimenti dai dati che cambiano nel tempo. Le applicazioni IoT e DevOps la telemetria industriale possono utilizzare database di serie temporali.
Colonna ampia	Amazon Keyspaces (per Apache Cassandra)	Utilizza tabelle, righe e colonne, ma a differenza di un database relazionale, i nomi e il formato delle colonne possono variare da riga a riga all'interno della stessa tabella. In genere, gli store colonnari sono utilizzati nelle applicazioni industriali su larga scala per la manutenzione delle apparecchiature, la gestione delle flotte e l'ottimizzazione dei percorsi.
Di libri mastri	Database Amazon Quantum Ledger (Amazon) QLDB	Fornisce un'autorità centralizzata e affidabile per mantenere un registro delle transazioni scalabile, immutabile e verificabile tramite crittografia per ogni applicazione. I database di libri mastri vengono utilizzati per sistemi di record, catena di fornitura, registrazioni e persino transazioni bancarie.

- Se stai creando una piattaforma dati, sfrutta un'[architettura di dati moderna](#) AWS per integrare il tuo data lake, il data warehouse e gli archivi dati creati appositamente.
- Le domande chiave da porsi quando si sceglie un data store per il carico di lavoro sono le seguenti:

Domanda	Aspetti da considerare
Come sono strutturati i dati?	<ul style="list-style-type: none"> • Se i dati non sono strutturati, prendi in considerazione un object store come Amazon S3 o un database No SQL come Amazon DocumentDB • Per i dati chiave-valore, prendi in considerazione DynamoDB, Amazon (ElastiCache Redis) o Amazon MemoryDB OSS
Quale livello di integrità referenziale è richiesto?	<ul style="list-style-type: none"> • Per i vincoli di chiave esterna, i database relazionali come Amazon e RDS Aurora possono fornire questo livello di integrità. • In genere, all'interno di un SQL modello senza dati, i dati vengono denormalizzati in un unico documento o in una raccolta di documenti da recuperare in un'unica richiesta anziché riunirli tra più documenti o tabelle.
È richiesta la conformità ACID (atomicità, coerenza, isolamento, durabilità)?	<ul style="list-style-type: none"> • Se le ACID proprietà associate ai database relazionali sono obbligatorie, prendi in considerazione un database relazionale come Amazon e RDS Aurora. • Se è richiesta una forte coerenza per No SQL database, è possibile utilizzare letture fortemente coerenti con DynamoDB.
Come cambierà nel tempo l'archiviazione? In che modo questo avrà effetto sulla scalabilità?	<ul style="list-style-type: none"> • I database serverless come DynamoDB e Amazon Quantum Ledger Database (Amazon) verranno scalati dinamicamente. QLDB

Domanda	Aspetti da considerare
	<ul style="list-style-type: none"> • Per i database relazionali sono previsti limiti massimi per l'archiviazione allocata, al raggiungimento dei quali si rende spesso necessario partizionare orizzontalmente tali database tramite meccanismi quali la partizione.
<p>Qual è la proporzione di query in lettura rispetto alle quelle in scrittura? Il caching potrebbe probabilmente migliorare le prestazioni?</p>	<ul style="list-style-type: none"> • I carichi di lavoro impegnativi in lettura possono trarre vantaggio da un livello di caching, ad esempio DAX se il database è ElastiCache DynamoDB. • Le letture possono anche essere scaricate per leggere le repliche con database relazionali come Amazon. RDS
<p>L'archiviazione e la modifica (OLTP- Elaborazione delle transazioni online) o il recupero e il reporting (OLAP- Elaborazione analitica online) hanno una priorità più elevata?</p>	<ul style="list-style-type: none"> • Per l'elaborazione transazionale read-as-is ad alto rendimento, prendi in considerazione un database No come DynamoDB. SQL • Per modelli di lettura complessi e ad alta velocità (come join) con coerenza, usa Amazon. RDS • Per le query analitiche, prendi in considerazione un database a colonne come Amazon Redshift o l'esportazione dei dati su Amazon S3 e l'esecuzione di analisi utilizzando Athena o Amazon. QuickSight

Domanda	Aspetti da considerare
Che livello di durabilità è necessario per i dati?	<ul style="list-style-type: none">• Aurora replica automaticamente i dati su tre zone di disponibilità all'interno di una regione, il che significa che i dati sono altamente durevoli con minori probabilità di perdite.• DynamoDB viene automaticamente replicato in più zone di disponibilità per offrire livelli elevati di disponibilità e durabilità dei dati.• Amazon S3 offre il 99,999999999 di durabilità. Molti servizi di database, come Amazon RDS e DynamoDB, supportano l'esportazione di dati su Amazon S3 per la conservazione e l'archiviazione a lungo termine.
È presente il desiderio di abbandonare i motori di database commerciali o i costi di licenza?	<ul style="list-style-type: none">• Prendi in considerazione motori open source come PostgreSQL e MySQL Amazon o RDS Aurora.• Sfrutta AWS Database Migration Service e AWS Schema Conversion Tool per eseguire le migrazioni dai motori di database commerciali a quelli open-source
Quali sono le aspettative operative per il database? Il passaggio ai servizi gestiti è una priorità?	<ul style="list-style-type: none">• Sfruttare Amazon RDS anziché Amazon EC2 e DynamoDB o Amazon DocumentDB anziché ospitare autonomamente SQL un database. Non può ridurre il sovraccarico operativo.

Domanda	Aspetti da considerare
<p>Come avviene attualmente l'accesso al database? Si tratta solo di accesso alle applicazioni o ci sono utenti di business intelligence (BI) e altre applicazioni connesse? off-the-shelf</p>	<ul style="list-style-type: none"> In presenza di dipendenze da strumenti esterni, potresti dover mantenere la compatibilità con i database che supportano. Amazon RDS è completamente compatibile con le diverse versioni del motore che supporta, tra cui Microsoft SQL Server, OracleSQL, My e SQL Postgre.

- Esegui esperimenti e benchmarking in un ambiente non di produzione per identificare quale datastore può soddisfare al meglio i requisiti del tuo carico di lavoro.

Risorse

Documenti correlati:

- [Tipi di EBS volume Amazon](#)
- [EC2Archiviazione Amazon](#)
- [AmazonEFS: EFS Prestazioni Amazon](#)
- [Prestazioni FSx di Amazon for Lustre](#)
- [Prestazioni FSx di Amazon per Windows File Server](#)
- [Amazon S3 Glacier: documentazione di S3 Glacier](#)
- [Amazon S3: considerazioni su velocità e prestazioni delle richieste](#)
- [Archiviazione su cloud con AWS](#)
- [Caratteristiche di Amazon EBS I/O](#)
- [Database su cloud con AWS](#)
- [AWS Database Caching](#)
- [DynamoDB Accelerator](#)
- [Best practice di Amazon Aurora](#)
- [Prestazioni di Amazon RedShift](#)
- [Amazon Athena top 10 performance tips](#)
- [Amazon Redshift Spectrum best practices](#)

- [Amazon DynamoDB best practices](#)
- [Scegli tra Amazon EC2 e Amazon RDS](#)
- [Best practice per l'implementazione di Amazon ElastiCache](#)

Video correlati:

- [AWS re:Invent 2023: migliora l'efficienza di Amazon Elastic Block Store e sii più efficiente in termini di costi](#)
- [AWS re:Invent 2023: ottimizzazione del prezzo e delle prestazioni dello storage con Amazon Simple Storage Service](#)
- [AWS re:Invent 2023: creazione e ottimizzazione di un data lake su Amazon Simple Storage Service](#)
- [AWS re:Invent 2022: creazione di moderne architetture di dati su AWS](#)
- [AWS re:Invent 2022: creazione di architetture di data mesh su AWS](#)
- [AWS re:Invent 2023: approfondimenti su Amazon Aurora e sulle sue innovazioni](#)
- [AWS re:Invent 2023: modellazione avanzata dei dati con Amazon DynamoDB](#)
- [AWS re:Invent 2022: modernizza le app con database creati appositamente](#)
- [Amazon DynamoDB deep dive: Advanced design patterns](#)

Esempi correlati:

- [AWS Workshop sui database creati appositamente](#)
- [Databases for Developers](#)
- [AWS Giornata di immersione nell'architettura dei dati moderna](#)
- [Crea una rete di dati su AWS](#)
- [Esempi di Amazon S3](#)
- [Optimize Data Pattern using Amazon Redshift Data Sharing](#)
- [Migrazioni dei database](#)
- [MS SQL Server - AWS Database Migration Service \(AWS DMS\) Demo sulla replica](#)
- [Database Modernization Hands On Workshop](#)
- [Esempi di Amazon Neptune](#)

PERF03-BP02 Valuta le opzioni di configurazione disponibili per l'archivio dati

Comprendi e valuta le varie funzionalità e opzioni di configurazione disponibili per i tuoi datastore per ottimizzare lo spazio di archiviazione e le prestazioni per il tuo carico di lavoro.

Anti-pattern comuni:

- Utilizzi solo un tipo di storage, ad esempio AmazonEBS, per tutti i carichi di lavoro.
- Utilizzi il provisioning IOPS per tutti i carichi di lavoro senza eseguire test reali su tutti i livelli di storage.
- Non conosci le opzioni di configurazione della soluzione di gestione dei dati scelta.
- Ti basi soltanto sull'aumento delle dimensioni dell'istanza, senza tenere conto di altre opzioni di configurazione disponibili.
- Non esegui il test delle caratteristiche di dimensionamento del tuo datastore.

Vantaggi dell'adozione di questa best practice: esplorare le configurazioni del datastore e sperimentare con esse può consentire di ridurre il costo dell'infrastruttura, migliorare le prestazioni e ridurre l'impegno richiesto per mantenere i carichi di lavoro.

Livello di rischio associato se questa best practice non fosse adottata: medio

Guida all'implementazione

Un carico di lavoro può utilizzare uno o più datastore in base ai requisiti di archiviazione di dati e relativo accesso. Per ottimizzare prestazioni, efficienza e costi, è necessario valutare gli schemi di accesso ai dati per determinare le configurazioni appropriate del datastore. Nella valutazione delle opzioni di datastore, prendi in considerazione vari aspetti come le opzioni di archiviazione, la memoria, l'elaborazione, la replica di lettura, i requisiti di coerenza, il pool di connessioni e le opzioni di caching. Esegui esperimenti con queste diverse opzioni di configurazione per migliorare i parametri di efficienza delle prestazioni.

Passaggi dell'implementazione

- Esamina le configurazioni correnti (come il tipo di istanza, la dimensione di archiviazione o la versione del motore di database) del tuo datastore.

- AWS Consulta la documentazione e le best practice per conoscere le opzioni di configurazione consigliate che possono contribuire a migliorare le prestazioni del tuo data store. Le principali opzioni da considerare per il datastore sono le seguenti:

Opzione di configurazione	Esempi
Riduzione del carico delle letture (come le repliche di lettura e la memorizzazione nella cache)	<ul style="list-style-type: none">• Per le tabelle DynamoDB, è possibile scaricare le letture utilizzando per la memorizzazione nella cache. DAX• Puoi creare un cluster Amazon ElastiCache (RedisOSS) e configurare l'applicazione in modo che legga prima dalla cache, per poi tornare al database se l'elemento richiesto non è presente.• I database relazionali come Amazon RDS e Aurora e i database forniti SQL No come Neptune e Amazon DocumentDB supportano tutti l'aggiunta di repliche di lettura per scaricare le parti di lettura del carico di lavoro.• I database serverless come DynamoDB si dimensionano automaticamente. Assicurati di disporre di unità di capacità di lettura sufficienti () per gestire il carico di lavoro. RCU

Opzione di configurazione	Esempi
Dimensionamento delle scritture (come la partizione delle chiavi di partizione o l'introduzione di una coda)	<ul style="list-style-type: none">• Per i database relazionali, è possibile aumentare le dimensioni dell'istanza per far fronte a un maggiore carico di lavoro o aumentare il provisioning IOPs per consentire una maggiore velocità di trasmissione dello storage sottostante.• È anche possibile introdurre una coda davanti al database, invece di eseguire direttamente la scrittura su di esso. Questo schema consente di disaccoppiare l'acquisizione dal database e controllare il flusso, in modo che il database sia in grado di gestirlo.• Raggruppare in batch le richieste di scrittura, anziché creare molte transazioni di breve durata, può aiutare a migliorare il throughput in database relazionali con un elevato volume in scrittura.• I database serverless come DynamoDB possono scalare il throughput di scrittura automaticamente o regolando le unità WCU di capacità di scrittura assegnate () a seconda della modalità di capacità.• È tuttavia possibile che si verifichino problemi con le partizioni hot quando si raggiungono i limiti di throughput per una determinata chiave di partizione. Questo problema può essere arginato scegliendo una chiave di partizione con una distribuzione più uniforme o eseguendo lo sharding in lettura della chiave di partizione.

Opzione di configurazione	Esempi
Policy per gestire il ciclo di vita dei set di dati	<ul style="list-style-type: none">• Puoi utilizzare il ciclo di vita Amazon S3 per gestire gli oggetti durante il loro ciclo di vita. In caso di schemi di accesso sconosciuti, mutevoli o imprevedibili, puoi utilizzare il Piano intelligente Amazon S3, che monitora gli schemi di accesso e sposta in automatico gli oggetti che non hanno fatto registrar e accessi a livelli di accessi più economici . Sfrutta i parametri di Amazon S3 Storage Lens per individuare opportunità di ottimizzazione e lacune nella gestione del ciclo di vita.• Amazon EFS Lifecycle Management gestisce automaticamente lo storage dei file per i tuoi file system.
Gestione e pooling delle connessioni	<ul style="list-style-type: none">• Amazon RDS Proxy può essere utilizzato con Amazon RDS e Aurora per gestire le connessioni al database.• I database serverless come DynamoDB non hanno connessioni associate, ma valuta la capacità assegnata e le policy di dimensionamento automatico per affrontare i picchi nel carico.

- Esegui esperimenti e benchmarking in un ambiente non di produzione per identificare quale opzione di configurazione può soddisfare i requisiti del tuo carico di lavoro.
- Dopo gli esperimenti, pianifica la migrazione e convalida i parametri delle prestazioni.
- Utilizza strumenti di AWS monitoraggio (come [Amazon CloudWatch](#)) e ottimizzazione (come [Amazon S3 Storage Lens](#)) per ottimizzare continuamente il tuo archivio dati utilizzando modelli di utilizzo reali.

Risorse

Documenti correlati:

- [Archiviazione nel cloud in AWS](#)
- [Tipi di EBS volume Amazon](#)
- [EC2Archiviazione Amazon](#)
- [AmazonEFS: EFS prestazioni Amazon](#)
- [Prestazioni Amazon FSx for Lustre](#)
- [Prestazioni FSx di Amazon per Windows File Server](#)
- [Amazon S3 Glacier: documentazione di S3 Glacier](#)
- [Amazon S3: considerazioni su velocità e prestazioni delle richieste](#)
- [Caratteristiche di Amazon EBS I/O](#)
- [Database su cloud con AWS](#)
- [AWS Database Caching](#)
- [DynamoDB Accelerator](#)
- [Best practice di Amazon Aurora](#)
- [Prestazioni di Amazon RedShift](#)
- [Amazon Athena top 10 performance tips](#)
- [Amazon Redshift Spectrum best practices](#)
- [Amazon DynamoDB best practices](#)

Video correlati:

- [AWS re:Invent 2023: Improve Amazon Elastic Block Store efficiency and be more cost-efficient](#)
- [AWS re:Invent 2023: Optimize storage price and performance with Amazon Simple Storage Service](#)
- [AWS re:Invent 2023: Building and optimizing a data lake on Amazon Simple Storage Service](#)
- [AWS re:Invent 2023: Cosa c'è di nuovo con lo storage di file AWS](#)
- [AWS re:Invent 2023: Dive deep into Amazon DynamoDB](#)

Esempi correlati:

- [AWS Workshop sui database appositamente progettati](#)
- [Databases for Developers](#)
- [AWS Giornata di immersione nell'architettura dei dati moderna](#)
- [Amazon EBS Scalabilità automatica](#)
- [Esempi di Amazon S3](#)
- [Esempi di Amazon DynamoDB](#)
- [AWS Esempi di migrazione di database](#)
- [Workshop sulla modernizzazione dei database](#)
- [Utilizzo dei parametri sul database Amazon RDS for Postgress](#)

PERF03-BP03 Raccolta e registrazione delle metriche delle prestazioni degli archivi dati

Tieni traccia e registra i parametri delle prestazioni pertinenti per il tuo datastore per capire l'andamento delle prestazioni delle soluzioni di gestione dei dati. Questi parametri possono aiutarti a ottimizzare il tuo datastore, verificare che i requisiti del carico di lavoro siano rispettati e fornire una panoramica chiara sull'andamento delle prestazioni del carico di lavoro.

Anti-pattern comuni:

- Utilizzi solo i file di log manuali per la ricerca dei parametri.
- Pubblichiamo i parametri solo sugli strumenti interni utilizzati dal tuo team e non hai un quadro completo del carico di lavoro.
- Utilizzo solo dei parametri predefiniti registrati dal software di monitoraggio selezionato.
- Revisione dei parametri solo quando c'è un problema.
- Monitori solo i parametri a livello di sistema, senza acquisire i parametri di accesso ai dati o di utilizzo.

Vantaggi dell'adozione di questa best practice: la definizione di una linea di base delle prestazioni ti aiuta a comprendere il comportamento normale e i requisiti dei carichi di lavoro. Gli schemi anomali possono essere identificati ed eliminati più rapidamente, per migliorare le prestazioni e l'affidabilità del datastore.

Livello di rischio associato se questa best practice non fosse adottata: elevato

Guida all'implementazione

Per monitorare le prestazioni dei datastore, devi registrare più parametri delle prestazioni in un periodo di tempo. Ciò consente di rilevare le anomalie e di misurare le prestazioni rispetto ai parametri aziendali, per verificare che le esigenze del carico di lavoro siano rispettate.

I parametri devono includere sia il sistema sottostante che supporta il datastore sia i parametri del database. Le metriche di sistema sottostanti possono includere CPU utilizzo, memoria, spazio di archiviazione su disco disponibile, I/O su disco, rapporto di accesso alla cache e metriche di rete in entrata e in uscita, mentre le metriche del data store possono includere transazioni al secondo, query principali, tassi medi di query, tempi di risposta, utilizzo dell'indice, blocchi delle tabelle, timeout delle query e numero di connessioni aperte. Questi dati sono cruciali per capire l'andamento del carico di lavoro e come viene utilizzata la soluzione di gestione dei dati. Utilizza tali parametri come parte di un approccio basato sui dati per mettere a punto e ottimizzare le risorse del tuo carico di lavoro.

Utilizza strumenti, librerie e sistemi che registrano misure delle prestazioni relative alle prestazioni del database.

Passaggi dell'implementazione

- Determina i principali parametri delle prestazioni da monitorare per il tuo datastore.
 - [Parametri e dimensioni di Amazon S3](#)
 - [Parametri di monitoraggio per un'istanza Amazon RDS](#)
 - [Monitoraggio del carico del DB con Performance Insights su Amazon RDS](#)
 - [Panoramica sul monitoraggio avanzato](#)
 - [DynamoDB Metrics and dimensions](#)
 - [Monitoraggio di DynamoDB Accelerator](#)
 - [Monitoraggio di Amazon MemoryDB con Amazon CloudWatch](#)
 - [Quali parametri è opportuno monitorare?](#)
 - [Monitoring Amazon Redshift cluster performance](#)
 - [Timestream metrics and dimensions](#)
 - [CloudWatch Parametri Amazon per Amazon Aurora](#)
 - [Creazione di log e monitoraggio in Amazon Keyspaces \(per Apache Cassandra\)](#)
 - [Monitoring Amazon Neptune Resources](#)

- Utilizza una soluzione di registrazione e monitoraggio approvata per raccogliere queste metriche. [Amazon CloudWatch](#) può raccogliere metriche tra le risorse della tua architettura. Puoi anche raccogliere e pubblicare parametri personalizzati per ottenere parametri aziendali o derivati. Utilizza soluzioni CloudWatch di terze parti per impostare allarmi che indicano quando vengono superate le soglie.
- Verifica se il monitoraggio dei datastore può trarre vantaggio da una soluzione di machine learning che rileva le anomalie delle prestazioni.
 - [Amazon DevOps Guru per Amazon RDS](#) offre visibilità sui problemi di prestazioni e fornisce consigli per azioni correttive.
- Configura la conservazione dei dati nella soluzione di monitoraggio e registrazione per soddisfare i tuoi obiettivi operativi e di sicurezza.
 - [Conservazione dei dati predefinita per le metriche CloudWatch](#)
 - [Conservazione dei dati predefinita per i registri CloudWatch](#)

Risorse

Documenti correlati:

- [AWS Database Caching](#)
- [Amazon Athena top 10 performance tips](#)
- [Amazon Aurora best practices](#)
- [DynamoDB Accelerator](#)
- [Amazon DynamoDB best practices](#)
- [Amazon Redshift Spectrum best practices](#)
- [Prestazioni di Amazon RedShift](#)
- [Database cloud con AWS](#)
- [Amazon RDS Performance Insights](#)

Video correlati:

- [AWS re:Invent 2022 - Monitoraggio delle prestazioni con Amazon e RDS Aurora, con Autodesk](#)
- [Monitoraggio e ottimizzazione delle prestazioni del database con Amazon DevOps Guru per Amazon RDS](#)
- [AWS re:Invent 2023 - Cosa c'è di nuovo con lo storage di file AWS](#)

- [AWS re:Invent 2023 - Scopri di più su Amazon DynamoDB](#)
- [AWS re:Invent 2023 - Creazione e ottimizzazione di un data lake su Amazon S3](#)
- [AWS re:Invent 2023 - Cosa c'è di nuovo con l'archiviazione dei file AWS](#)
- [AWS re:Invent 2023 - Scopri di più su Amazon DynamoDB](#)
- [Best practice per il monitoraggio dei carichi di lavoro Redis su Amazon ElastiCache](#)

Esempi correlati:

- [Framework di raccolta dei parametri di ingestione del set di dati AWS](#)
- [Workshop sul RDS monitoraggio di Amazon](#)
- [AWS Workshop sui database appositamente progettato](#)

PERF03-BP04 Implementazione di strategie per migliorare le prestazioni delle query nel datastore

Implementa le strategie per ottimizzare i dati e migliorare le query sui dati in modo da consentire una maggiore scalabilità e prestazioni più efficienti per il tuo carico di lavoro.

Anti-pattern comuni:

- Non suddividi i dati in partizioni nel tuo datastore.
- I dati vengono archiviati in un solo formato di file nel tuo datastore.
- Non usi gli indici nel tuo datastore.

Vantaggi dell'adozione di questa best practice: l'ottimizzazione delle prestazioni dei dati e delle query si traduce in maggiore efficienza, costi inferiori e migliore esperienza utente.

Livello di rischio associato se questa best practice non fosse adottata: medio

Guida all'implementazione

L'ottimizzazione di dati e query è un aspetto critico dell'efficienza delle prestazioni in un datastore, poiché influisce sulle prestazioni e sulla reattività dell'intero carico di lavoro cloud. Le query non ottimizzate possono comportare un maggiore utilizzo delle risorse e rallentamenti, riducendo così l'efficienza complessiva di un datastore.

L'ottimizzazione dei dati include diverse tecniche per garantire prestazioni efficienti per l'archiviazione di dati e il relativo accesso. Ciò aiuta anche a migliorare le prestazioni delle query in un datastore. Le strategie chiave includono il partizionamento, la compressione e la denormalizzazione dei dati, che contribuiscono a ottimizzare i dati sia per l'archiviazione che per l'accesso.

Passaggi dell'implementazione

- Esamina e analizza le query sui dati critiche che vengono eseguite nel tuo datastore.
- Individua le query lente del tuo datastore e utilizza i piani di query per comprenderne lo stato attuale.
 - [Analisi del piano di query in Amazon Redshift](#)
 - [Using EXPLAIN and EXPLAIN ANALYZE in Athena](#)
- Implementa le strategie per migliorare le prestazioni delle query. Ecco alcune strategie chiave:
 - Utilizzo di un [formato di file colonnare](#) (come Parquet o ORC).
 - Compressione dei dati nel datastore per ridurre lo spazio di archiviazione e il funzionamento di I/O.
 - Partizionamento dei dati per suddividere i dati in parti più piccole e ridurre i tempi di analisi dei dati.
 - [Partizionamento dei dati in Athena](#)
 - [Partitions and data distribution](#)
 - Indicizzazione dei dati sulle colonne comuni della query.
 - Uso delle viste materializzate per le domande frequenti.
 - [Understanding materialized views](#)
 - [Creating materialized views in Amazon Redshift](#)
 - Scelta dell'operazione di unione corretta per la query. Quando unisci due tabelle, specifica la tabella più grande sul lato sinistro dell'unione e la tabella più piccola sul lato destro.
 - Miglioramento della latenza e riduzione del numero di operazioni di I/O del database grazie alla soluzione di cache distribuita.
 - Manutenzione regolare, ad esempio [vacuum](#), reindicizzazione ed [esecuzione di statistiche](#).
- La sperimentazione e i test delle strategie in un ambiente non di produzione.

Risorse

Documenti correlati:

- [Best practice di Amazon Aurora](#)
- [Prestazioni di Amazon RedShift](#)
- [Amazon Athena top 10 performance tips](#)
- [AWS Database Caching](#)
- [Best practice per l'implementazione di Amazon ElastiCache](#)
- [Partizionamento dei dati in Athena](#)

Video correlati:

- [AWS re:Invent 2023 - AWS storage cost-optimization best practices](#)
- [AWS re:Invent 2022 - Performance monitoring with Amazon RDS and Aurora, featuring Autodesk](#)
- [Optimize Amazon Athena Queries with New Query Analysis Tools](#)

Esempi correlati:

- [AWS Purpose Built Databases Workshop](#)

PERF03-BP05 Implementare modelli di accesso ai dati che utilizzano la memorizzazione nella cache

Implementa modelli di accesso che possano trarre vantaggio dalla memorizzazione dei dati nella cache per il recupero rapido dei dati a cui si accede di frequente.

Anti-pattern comuni:

- Memorizzare nella cache dati che cambiano in maniera frequente.
- Fare affidamento sui dati memorizzati nella cache come se fossero archiviati in modo duraturo e sempre disponibili.
- Non tenere conto della coerenza dei dati memorizzati nella cache.
- Non monitorare l'efficienza dell'implementazione della cache.

Vantaggi dell'adozione di questa best practice: l'archiviazione dei dati in una cache può migliorare la latenza di lettura, il throughput, l'esperienza utente e l'efficienza complessiva, oltre a ridurre i costi.

Livello di rischio associato se questa best practice non fosse adottata: medio

Guida all'implementazione

Una cache è un componente software o hardware progettato per archiviare dati in modo che le richieste future degli stessi dati possano essere soddisfatte più velocemente o in modo più efficiente. I dati memorizzati in una cache possono essere ricostruiti in caso di perdita, ripetendo un calcolo precedente o recuperandolo da un altro datastore.

La memorizzazione dei dati nella cache può essere una delle strategie più efficaci per migliorare le prestazioni complessive delle applicazioni e ridurre il carico sulle origini dati primarie sottostanti. È possibile memorizzare i dati nella cache a più livelli dell'applicazione, ad esempio all'interno dell'applicazione che effettua chiamate remote, operazione nota come memorizzazione nella cache lato client, o mediante un servizio secondario veloce per l'archiviazione dei dati, operazione nota come memorizzazione nella cache remota.

Memorizzazione nella cache lato client

Con la memorizzazione nella cache lato client, ogni client (un'applicazione o un servizio che interroga il datastore di backend) può archiviare localmente i risultati delle proprie query uniche per un periodo di tempo specificato. Ciò può ridurre il numero di richieste a un datastore attraverso la rete perché viene controllata prima la cache del client locale. Se questa non contiene risultati, l'applicazione può interrogare il datastore e archiviare tali risultati localmente. Questo modello consente a ciascun client di archiviare i dati nella sede più vicina possibile (il client stesso), garantendo così la latenza più bassa possibile. I client possono inoltre continuare a eseguire query quando il datastore di backend non è disponibile, aumentando la disponibilità dell'intero sistema.

Uno svantaggio di questo approccio è che quando sono coinvolti più client, potrebbero archiviare localmente gli stessi dati memorizzati nella cache. Ciò si traduce in un utilizzo duplicato dell'archiviazione e nell'incoerenza dei dati tra questi client. Può accadere che un client memorizzi nella cache i risultati di una query e un minuto dopo un altro client esegua la stessa query ottenendo un risultato diverso.

Memorizzazione nella cache remota

Come soluzione al problema della duplicazione dei dati tra client, è possibile utilizzare un servizio esterno veloce o la memorizzazione nella cache remota per archiviare i dati sottoposti a query. Aniché controllare un datastore locale, ogni client controllerà la cache remota prima di interrogare il datastore di backend. Questa strategia consente di ottenere risposte più coerenti tra i client, una

migliore efficienza dei dati archiviati e un volume maggiore di dati memorizzati nella cache, perché lo spazio di archiviazione si dimensiona in maniera indipendente dai client.

Lo svantaggio di una cache remota è che l'intero sistema può registrare una latenza più elevata, poiché è necessario un hop di rete aggiuntivo per controllare la cache remota. Per migliorare la latenza, è possibile utilizzare la memorizzazione nella cache lato client insieme alla memorizzazione nella cache remota, eseguendo così una memorizzazione nella cache su più livelli.

Passaggi dell'implementazione

- Identifica i database APIs e i servizi di rete che potrebbero trarre vantaggio dalla memorizzazione nella cache. I servizi che hanno carichi di lavoro di lettura elevati, hanno un read-to-write rapporto elevato o sono costosi da scalare sono candidati alla memorizzazione nella cache.
 - [Database Caching](#)
 - [Abilitare la API memorizzazione nella cache per migliorare la reattività](#)
- Identifica il tipo di strategia di memorizzazione nella cache più adatto al tuo modello di accesso.
 - [Caching strategies](#)
 - [AWS Caching Solutions](#)
- Attieniti alle [best practice sulla memorizzazione nella cache](#) per il tuo archivio dati.
- Configura una strategia di invalidazione della cache, ad esempio a time-to-live (TTL), per tutti i dati che bilanci l'aggiornamento dei dati e riduca la pressione sul datastore di backend.
- Abilita funzionalità quali tentativi di connessione automatici, backoff esponenziale, timeout lato client e pool di connessioni nel client, se disponibili, che possono migliorare prestazioni e affidabilità.
 - [Migliori pratiche: client Redis e Amazon ElastiCache \(OSSRedis\)](#)
- Monitora la percentuale di riscontri nella cache con un obiettivo dell'80% o superiore. Valori inferiori possono indicare una dimensione della cache insufficiente o un modello di accesso che non sfrutta la memorizzazione nella cache.
 - [Which metrics should I monitor?](#)
 - [Le migliori pratiche per il monitoraggio dei carichi di lavoro Redis su Amazon ElastiCache](#)
 - [Monitoraggio delle best practice con Amazon ElastiCache \(RedisOSS\) tramite Amazon CloudWatch](#)
- Implementa la [replica dei dati](#) per eliminare il carico delle letture per più istanze e migliorare prestazioni e disponibilità della lettura dei dati.

Risorse

Documenti correlati:

- [Utilizzo dell'obiettivo Amazon ElastiCache Well-Architected](#)
- [Monitoraggio delle best practice con Amazon ElastiCache \(RedisOSS\) tramite Amazon CloudWatch](#)
- [Quali parametri è opportuno monitorare?](#)
- [Prestazioni su larga scala con il ElastiCache white paper di Amazon](#)
- [Sfide e strategie del caching](#)

Video correlati:

- [Percorso di ElastiCache apprendimento su Amazon](#)
- [Progetta per il successo con le ElastiCache best practice di Amazon](#)
- [AWS re:Invent 2020 - Progetta per il successo con le best practice di Amazon ElastiCache](#)
- [AWS re:Invent 2023 - \[\] LAUNCH Presentazione di Amazon Serverless ElastiCache](#)
- [AWS re:Invent 2022 - 5 ottimi modi per reimmaginare il tuo livello di dati con Redis](#)
- [AWS re:Invent 2021 - Approfondimento su Amazon ElastiCache \(Redis\) OSS](#)

Esempi correlati:

- [Potenziamento delle prestazioni SQL del mio database con Amazon ElastiCache \(RedisOSS\)](#)

Reti e distribuzione di contenuti

La soluzione di rete ottimale per un carico di lavoro varia in base a latenza, requisiti di throughput, jitter e larghezza di banda. I vincoli fisici, ad esempio le risorse utente o on-premises, determinano le opzioni di posizione. Questi vincoli possono essere compensati con le posizioni edge o la collocazione delle risorse.

In AWS, le reti sono virtualizzate e vengono fornite in molti tipi e configurazioni diversi. In questo modo puoi soddisfare le tue esigenze di rete più facilmente. AWS offre funzionalità di prodotto (ad esempio reti avanzate, istanze Amazon EC2 ottimizzate per la rete, accelerazione del trasferimento Amazon S3 e Amazon CloudFront dinamico) pensate per l'ottimizzazione del traffico di rete. AWS offre anche funzionalità di rete (ad esempio instradamento in base alla latenza di Amazon Route 53, endpoint VPC Amazon, AWS Direct Connect e AWS Global Accelerator) per ridurre la distanza di rete o il jitter.

Questa area di interesse offre linee guida e best practice per progettare, configurare e gestire soluzioni di rete e distribuzione di contenuti nel cloud in maniera efficiente.

Best practice

- [PERF04-BP01 Scopri come la rete influisce sulle prestazioni](#)
- [PERF04-BP02 Valuta le funzionalità di rete disponibili](#)
- [PERF04-BP03 Scegli la connettività dedicata appropriata o per il tuo carico di lavoro VPN](#)
- [PERF04-BP04 Usa il bilanciamento del carico per distribuire il traffico su più risorse](#)
- [PERF04-BP05 Scegli i protocolli di rete per migliorare le prestazioni](#)
- [PERF04-BP06 Scegli la posizione del carico di lavoro in base ai requisiti di rete](#)
- [PERF04-BP07 Ottimizza la configurazione di rete in base a metriche](#)

PERF04-BP01 Scopri come la rete influisce sulle prestazioni

Analizza e comprendi in che modo le decisioni correlate alla rete influiscono sul carico di lavoro per fornire prestazioni efficienti e una migliore esperienza utente.

Anti-pattern comuni:

- Tutto il traffico passa attraverso i data center esistenti.

- Si instrada tutto il traffico attraverso i firewall centrali anziché utilizzare strumenti di sicurezza di rete nativi del cloud.
- Effettua il provisioning delle AWS Direct Connect connessioni senza comprendere i requisiti di utilizzo effettivi.
- Quando si definiscono le soluzioni di rete, non si considerano le caratteristiche del carico di lavoro e l'overhead della crittografia.
- Per le soluzioni di rete nel cloud si utilizzano concetti e strategie on-premises.

Vantaggi dell'adozione di questa best practice: la comprensione dell'impatto della rete sulle prestazioni del carico di lavoro ti aiuta a identificare i potenziali colli di bottiglia, migliorare l'esperienza dell'utente, aumentare l'affidabilità e ridurre la manutenzione operativa al variare del carico di lavoro.

Livello di rischio associato se questa best practice non fosse adottata: elevato

Guida all'implementazione

La rete è responsabile della connettività tra componenti dell'applicazione, servizi cloud, reti edge e dati on-premises e quindi può avere un forte impatto sulle prestazioni dei carichi di lavoro. Oltre alle prestazioni del carico di lavoro, l'esperienza dell'utente può essere influenzata anche da latenza della rete, larghezza di banda, protocolli, posizione, congestione della rete, jitter, throughput e regole di instradamento.

Predisponi un elenco documentato dei requisiti di rete del carico di lavoro, tra cui latenza, dimensione dei pacchetti, regole di instradamento, protocolli e modelli di traffico di supporto. Esamina le soluzioni di rete disponibili e individua il servizio che soddisfi le caratteristiche di rete del proprio carico di lavoro. Le reti basate sul cloud possono essere ricostruite rapidamente, quindi l'evoluzione dell'architettura di rete nel tempo è necessaria per migliorare l'efficienza delle prestazioni.

Passaggi dell'implementazione:

- Definisci e documenta i requisiti di prestazioni di rete, tra cui metriche come latenza di rete, larghezza di banda, protocolli, posizioni, modelli di traffico (picchi e frequenza), throughput, crittografia, ispezione e regole di instradamento.
- Scopri i principali servizi AWS di rete come [VPCs](#), [Elastic Load Balancing \(ELB\)](#) e [Amazon Route 53](#). [AWS Direct Connect](#)
- Acquisisci le seguenti caratteristiche di rete fondamentali:

Caratteristiche	Strumenti e metriche
Caratteristiche fondamentali della rete	<ul style="list-style-type: none"> • VPCRegistri di flusso • AWS Transit Gateway Log di flusso • AWS Transit Gateway metriche • AWS PrivateLink metriche
Caratteristiche di rete dell'applicazione	<ul style="list-style-type: none"> • Elastic Fabric Adapter (EFA) • AWS App Mesh metriche • Metriche di Amazon API Gateway
Caratteristiche della rete edge	<ul style="list-style-type: none"> • CloudFront Metriche Amazon • Parametri di Amazon Route 53 • AWS Global Accelerator metriche
Caratteristiche della rete ibrida	<ul style="list-style-type: none"> • AWS Direct Connect metriche • AWS Site-to-Site VPN metriche • AWS Client VPN metriche • Cloud AWS WANmetriche
Caratteristiche della sicurezza di rete	<ul style="list-style-type: none"> • AWS Shield e AWS WAF metriche AWS Network Firewall
Caratteristiche del tracciamento	<ul style="list-style-type: none"> • AWS X-Ray • VPCReachability Analyzer • Strumento di analisi degli accessi alla rete • Amazon Inspector • Amazon CloudWatch RUM

- Esegui il benchmark e testa le prestazioni della rete:
 - Effettua un [benchmark](#) del throughput di rete, poiché alcuni fattori possono influire sulle prestazioni EC2 della rete Amazon quando le istanze sono uguali. VPC Misura la larghezza di banda di rete tra le istanze di Amazon EC2 Linux nella stessa istanza. VPC
 - Effettua [test di carico](#) per sperimentare soluzioni e opzioni di rete.

Risorse

Documenti correlati:

- [Application Load Balancer](#)
- [EC2Rete avanzata su Linux](#)
- [EC2Rete avanzata su Windows](#)
- [EC2Gruppi di collocamento](#)
- [Abilitazione di reti avanzate con Elastic Network Adapter \(ENA\) su istanze Linux](#)
- [Network Load Balancer](#)
- [Prodotti di rete con AWS](#)
- [Gateway di transito](#)
- [Transitioning to latency-based routing in Amazon Route 53](#)
- [VPCEndpoint](#)

Video correlati:

- [AWS re:Invent 2023 - fondamenti per il networking AWS](#)
- [AWS re:Invent 2023 - Cosa può fare il networking per la tua applicazione?](#)
- [AWS re:Invent 2023 - Design avanzati e nuove funzionalità VPC](#)
- [AWS re:Invent 2023 - Una guida per sviluppatori al cloud networking](#)
- [AWS re:Invent 2019 - Connettività e architetture di rete ibride AWSAWS](#)
- [AWS re:Invent 2019 - Ottimizzazione delle prestazioni di rete per le istanze Amazon EC2](#)
- [AWS Summit Online - Migliora le prestazioni della rete globale per le applicazioni](#)
- [AWS re:Invent 2020 - Migliori pratiche e suggerimenti per il networking con il Well-Architected Framework](#)
- [AWS re:Invent 2020 - migliori pratiche di rete nelle migrazioni su larga scala AWS](#)

Esempi correlati:

- [AWS Transit Gateway e soluzioni di sicurezza scalabili](#)
- [AWS Workshop di networking](#)
- [Hands-on Network Firewall Workshop](#)

- [Osservazione e diagnosi della rete su AWS](#)
- [Individuazione e risoluzione degli errori di configurazione della rete su AWS](#)

PERF04-BP02 Valuta le funzionalità di rete disponibili

Valuta le funzionalità di rete nel cloud che possono aumentare le prestazioni. Misura l'impatto di tali funzionalità attraverso test, parametri e analisi. Ad esempio, sfrutta le funzionalità a livello di rete disponibili per ridurre latenza, distanza di rete o jitter.

Anti-pattern comuni:

- Rimani all'interno di una regione perché è lì che si trova fisicamente la tua sede centrale.
- Utilizzi i firewall anziché i gruppi di sicurezza per filtrare il traffico.
- Fai una pausa TLS per ispezionare il traffico anziché affidarti a gruppi di sicurezza, policy degli endpoint e altre funzionalità native del cloud.
- Utilizzi solo la segmentazione basata su sottoreti anziché i gruppi di sicurezza.

Vantaggi dell'adozione di questa best practice la valutazione di tutte le funzionalità e le opzioni del servizio consente di ridurre il costo dell'infrastruttura e l'impegno necessario per mantenere il carico di lavoro e aumentare l'assetto di sicurezza generale. Puoi utilizzare il AWS backbone globale per fornire un'esperienza di rete ottimale ai tuoi clienti.

Livello di rischio associato se questa best practice non fosse adottata: elevato

Guida all'implementazione

AWS offre servizi come [AWS Global Accelerator](#) [Amazon CloudFront](#) che possono aiutare a migliorare le prestazioni di rete, mentre la maggior parte dei AWS servizi dispone di caratteristiche di prodotto (come la funzionalità [Amazon S3 Transfer Acceleration](#)) per ottimizzare il traffico di rete.

Analizza quali opzioni di configurazione relative alla rete sono disponibili e come possono influire sul tuo carico di lavoro. L'ottimizzazione delle prestazioni dipende dalla comprensione del modo in cui queste opzioni interagiscono con l'architettura e dall'impatto che hanno sulle prestazioni misurate e sull'esperienza utente.

Passaggi dell'implementazione

- Crea l'elenco dei componenti del carico di lavoro.

- Prendi in considerazione l'idea [Cloud AWS WAN](#) di utilizzarla per creare, gestire e monitorare la rete della tua organizzazione quando crei una rete globale unificata.
- Monitora le tue reti globali e principali con i [parametri di Amazon CloudWatch Logs](#). Sfrutta [Amazon CloudWatch RUM](#), che fornisce informazioni utili per identificare, comprendere e migliorare l'esperienza digitale degli utenti.
- Visualizza la latenza di rete aggregata tra Regioni AWS e le zone di disponibilità, nonché all'interno di ciascuna zona di disponibilità, utilizzala [AWS Network Manager](#) per ottenere informazioni dettagliate su come le prestazioni delle applicazioni si relazionano con le prestazioni della rete sottostante. AWS
- Utilizzate uno strumento esistente per il database di gestione della configurazione (CMDB) o un servizio, [AWS Config](#) ad esempio, per creare un inventario del carico di lavoro e della sua configurazione.
- Se si tratta di un carico di lavoro esistente, individua e documenta l'analisi di benchmark per le metriche relative alle prestazioni, concentrandoti sui colli di bottiglia e sulle aree da migliorare. Le metriche relative alla rete a livello di prestazioni varieranno a seconda dei requisiti aziendali e delle caratteristiche del carico di lavoro. Come punto di partenza, le seguenti metriche possono essere importanti per la revisione del carico di lavoro: larghezza di banda, latenza, perdita di pacchetti, jitter e ritrasmissioni.
- Se si tratta di un nuovo carico di lavoro, esegui [test di carico](#) per individuare i colli di bottiglia delle prestazioni.
- Per tutti i colli di bottiglia di questo tipo individuati, esamina le opzioni di configurazione per le soluzioni in uso per individuare le opportunità di miglioramento delle prestazioni. Consulta le seguenti opzioni e funzionalità di rete fondamentali:

Opportunità di miglioramento	Soluzione
Percorso o instradamenti di rete	Usa lo Strumento di analisi degli accessi alla rete per identificare percorsi o percorsi.
Protocolli di rete	Consulta la sezione PERF04-BP05 Scegli i protocolli di rete per migliorare le prestazioni
Topologia di rete	Valuta i compromessi operativi e prestazionali tra VPC Peering e AWS Transit Gateway quando si connettono più account. AWS Transit Gateway semplifica il modo in

Opportunità di miglioramento	Soluzione
	<p>cui interconnetti tutti i tuoi dispositivi VPCs, che possono estendersi su migliaia di reti locali. Account AWS Condividi il tuo AWS Transit Gateway tra più account utilizzando. AWS Resource Access Manager</p> <p>Consulta la sezione PERF04-BP03 Scegli la connettività dedicata appropriata o per il tuo carico di lavoro VPN</p>

Opportunità di miglioramento	Soluzione
Servizi di rete	<p>AWS Global Accelerator è un servizio di rete che migliora le prestazioni del traffico degli utenti fino al 60% utilizzando l'infrastruttura di rete AWS globale.</p> <p>Amazon CloudFront può migliorare le prestazioni della distribuzione e della latenza dei contenuti del carico di lavoro a livello globale.</p> <p>Usa Lambda @edge per eseguire funzioni che personalizzano i contenuti per renderli più vicini agli utenti, ridurre la latenza e migliorare le CloudFront prestazioni.</p> <p>Amazon Route 53 offre opzioni di instradamento basato sulla latenza, instradamento basato sulla geolocalizzazione, instradamento basato sulla geoprossimità e instradamento basato su IP per migliorare le prestazioni del tuo carico di lavoro per un pubblico globale. Rivedi il traffico del carico di lavoro e la posizione dell'utente quando il carico di lavoro è distribuito a livello globale per individuare quale opzione di instradamento è in grado di ottimizzare le prestazioni del carico di lavoro.</p>

Opportunità di miglioramento	Soluzione
Funzionalità delle risorse di archiviazione	<p>Amazon S3 Transfer Acceleration è una funzionalità che consente agli utenti esterni di beneficiare CloudFront delle ottimizzazioni di rete per caricare dati su Amazon S3. Ciò migliora le caratteristiche di trasferimento di grandi quantità di dati da posizioni remote prive di connettività dedicata al Cloud AWS.</p> <p>I punti di accesso multi-regione di Amazon S3 rappresentano una funzionalità che replica i contenuti in più regioni e semplifica il carico di lavoro fornendo un punto di accesso. Quando viene utilizzato un punto di accesso multi-regione, puoi richiedere o scrivere dati in Amazon S3 con il servizio che identifica il bucket con latenza più bassa.</p>

Opportunità di miglioramento	Soluzione
Funzionalità delle risorse di calcolo	<p>Le interfacce di rete elastiche (ENA) utilizzate dalle EC2 istanze, dai contenitori e dalle funzioni Lambda di Amazon sono limitate in base al flusso. Controlla i tuoi gruppi di collocamento per ottimizzare la velocità di trasmissione della rete. EC2 Per evitare colli di bottiglia a livello di flusso, progetta l'applicazione in modo che utilizzi più flussi. Per monitorare e ottenere visibilità sulle metriche di rete relative all'elaborazione, usa CloudWatch Metrics ed ethtool. Il ethtool comando è incluso nel ENA driver ed espone metriche aggiuntive relative alla rete che possono essere pubblicate come metriche personalizzate su. CloudWatch</p> <p>Amazon Elastic Network Adapters (ENA) forniscono un'ulteriore ottimizzazione offrendo un throughput migliore per le istanze all'interno di un gruppo di collocamento di cluster.</p> <p>Elastic Fabric Adapter (EFA) è un'interfaccia di rete per EC2 istanze Amazon che consente di eseguire carichi di lavoro che richiedono alti livelli di comunicazioni tra nodi su larga scala. AWS</p> <p>Le istanze EBS ottimizzate per Amazon utilizzano uno stack di configurazione ottimizzato e forniscono capacità aggiuntiva dedicata per aumentare l'I/O di Amazon. EBS</p>

Risorse

Documenti correlati:

- [Application Load Balancer](#)
- [EC2Rete avanzata su Linux](#)
- [EC2Rete avanzata su Windows](#)
- [EC2Gruppi di collocamento](#)
- [Abilitazione di reti avanzate con Elastic Network Adapter \(ENA\) su istanze Linux](#)
- [Network Load Balancer](#)
- [Prodotti di rete con AWS](#)
- [Transitioning to Latency-Based Routing in Amazon Route 53](#)
- [VPCEndpoint](#)
- [Log di flusso VPC](#)

Video correlati:

- [AWS re:Invent 2023 — Pronti per il futuro? Designing networks for growth and flexibility](#)
- [AWS re:Invent 2023 — Design avanzati e nuove funzionalità VPC](#)
- [AWS re:Invent 2023: guida per sviluppatori al cloud networking](#)
- [AWS re:Invent 2022 — Approfondisci l'infrastruttura di rete AWS](#)
- [AWS re:Invent 2019 — Connettività e architetture di rete ibride AWSAWS](#)
- [AWS re:Invent 2018 — Ottimizzazione delle prestazioni di rete per le istanze Amazon EC2](#)
- [AWS Global Accelerator](#)

Esempi correlati:

- [AWS Transit Gateway e soluzioni di sicurezza scalabili](#)
- [AWS Workshop di networking](#)
- [Observing and diagnosing your network](#)
- [Individuazione e risoluzione degli errori di configurazione della rete su AWS](#)

PERF04-BP03 Scegli la connettività dedicata appropriata o per il tuo carico di lavoro VPN

Quando hai bisogno di una connettività ibrida per connettere risorse on-premises e cloud, assicurati di avere una larghezza di banda adeguata per soddisfare i tuoi requisiti di prestazione. Fai una stima dei requisiti di larghezza di banda e latenza per il carico di lavoro ibrido. I valori calcolati determineranno le tue esigenze di dimensionamento.

Anti-pattern comuni:

- Valutate solo VPN soluzioni per i vostri requisiti di crittografia di rete.
- Non vengono valutate opzioni di backup o di connettività ridondante.
- Non è possibile identificare tutti i requisiti del carico di lavoro (esigenze di crittografia, protocollo, larghezza di banda e traffico).

Vantaggi dell'adozione di questa best practice: la selezione e la configurazione di soluzioni di connettività appropriate migliorano l'affidabilità del carico di lavoro e massimizzano le prestazioni. Identificando i requisiti del carico di lavoro, pianificando in anticipo e valutando le soluzioni ibride, è possibile ridurre al minimo le costose modifiche alla rete fisica e il sovraccarico operativo, aumentando al contempo i costi. time-to-value

Livello di rischio associato se questa best practice non fosse adottata: elevato

Guida all'implementazione

Sviluppa un'architettura di rete ibrida basata sui tuoi requisiti di larghezza di banda. [AWS Direct Connect](#) consente di connettere la rete on-premises in privato con AWS. È utile quando hai bisogno di larghezza di banda elevata, bassa latenza e di mantenere le prestazioni coerenti. Una VPN connessione stabilisce una connessione sicura su Internet. Viene utilizzata quando è necessaria solo una connessione temporanea, quando il costo è un fattore importante o come misura di contingenza in attesa che venga stabilita una connettività di rete fisica resiliente mentre AWS Direct Connect è in uso.

Se i requisiti di larghezza di banda sono elevati, potresti prendere in considerazione più servizi AWS Direct Connect o VPN più servizi. È possibile bilanciare il carico del traffico tra i servizi, anche se non consigliamo il bilanciamento del carico tra AWS Direct Connect e a VPN causa delle differenze di latenza e larghezza di banda.

Passaggi dell'implementazione

- Calcola i requisiti di larghezza di banda e latenza delle tue app esistenti.
 - Per i carichi di lavoro esistenti che stanno per essere trasferiti AWS, sfrutta i dati dei tuoi sistemi di monitoraggio della rete interna.
 - Per i carichi di lavoro nuovi o esistenti per i quali non sono disponibili dati di monitoraggio, consulta i proprietari dei prodotti per definire metriche sulle prestazioni adeguate e offrire un'esperienza utente soddisfacente.
- Seleziona una connessione dedicata o VPN come opzione di connettività. In base a tutti i requisiti del carico di lavoro (crittografia, larghezza di banda e traffico), puoi scegliere AWS Direct Connect o [AWS VPN](#) (o entrambi). Il diagramma seguente può aiutarti a scegliere il tipo di connessione appropriato.
 - [AWS Direct Connect](#) fornisce connettività dedicata all'ambiente AWS da 50 Mbps fino a 100 Gbps, utilizzando connessioni dedicate od ospitate. In questo modo, disporrai di latenza gestita e controllata, nonché di larghezza di banda assegnata, in modo che il carico di lavoro possa connettersi con efficienza ad altri ambienti. Utilizzando AWS Direct Connect i partner, è possibile disporre di end-to-end connettività da più ambienti, fornendo una rete estesa con prestazioni costanti. AWS offre la scalabilità della larghezza di banda della connessione Direct Connect utilizzando 100 Gbps nativi, link aggregation group (LAG) o BGP equal-cost multipath (). ECMP
 - AWS [Site-to-Site VPN](#) Fornisce un servizio gestito VPN che supporta la sicurezza del protocollo Internet (). IPsec Quando viene creata una VPN connessione, ogni VPN connessione include due tunnel per un'elevata disponibilità.
- Segui AWS la documentazione per scegliere l'opzione di connettività appropriata:
 - Se decidi di utilizzarla AWS Direct Connect, seleziona la larghezza di banda appropriata per la tua connettività.
 - Se utilizzi una connessione AWS Site-to-Site VPN tra più postazioni per connetterti a una Regione AWS, utilizza una [Site-to-SiteVPNconnessione accelerata](#) per avere l'opportunità di migliorare le prestazioni della rete.
 - Se la progettazione della rete prevede una IPsec VPN connessione via cavo [AWS Direct Connect](#), prendete in considerazione l'utilizzo di Private IP VPN per migliorare la sicurezza e ottenere la segmentazione. AWS Site-to-Site VPN L'[IP privato](#) viene distribuito sopra l'interfaccia virtuale di transito (VIF).
 - [AWS Direct Connect SiteLink](#) consente di creare connessioni ridondanti e a bassa latenza tra i data center in tutto il mondo inviando i dati lungo il percorso più veloce tra le sedi, evitando così di farlo. [AWS Direct Connect](#) Regioni AWS

- Convalida la configurazione della connettività prima di eseguire l'implementazione in produzione. Esegui test di sicurezza e prestazioni per assicurarti di soddisfare i requisiti di larghezza di banda, affidabilità, latenza e conformità.
- Monitora regolarmente le prestazioni e l'utilizzo della connettività e ottimizzali, se necessario.

Diagramma di flusso per le prestazioni deterministiche

Risorse

Documenti correlati:

- [Prodotti di rete con AWS](#)
- [AWS Transit Gateway](#)
- [VPC Endpoint](#)
- [Creazione di un'infrastruttura multirete scalabile e sicura VPC AWS](#)
- [Cliente VPN](#)

Video correlati:

- [AWS re:Invent 2023 — Creazione di una connettività di rete ibrida con AWS](#)
- [AWS re:Invent 2023 — Connettività remota sicura a AWS](#)
- [AWS re:Invent 2022 — Ottimizzazione delle prestazioni con Amazon CloudFront](#)
- [AWS re:Invent 2019 — Connettività e architetture di rete ibride AWS/AWS](#)
- [AWS re:Invent 2020 — Connect AWS Transit Gateway](#)

Esempi correlati:

- [AWS Transit Gateway e soluzioni di sicurezza scalabili](#)
- [AWS Workshop di networking](#)

PERF04-BP04 Usa il bilanciamento del carico per distribuire il traffico su più risorse

Distribuisce il traffico tra varie risorse o servizi affinché il carico di lavoro possa trarre vantaggio dall'elasticità fornita dal cloud. Puoi anche utilizzare il bilanciamento del carico per la terminazione dell'offloading della crittografia al fine di migliorare le prestazioni, l'affidabilità e gestire e instradare il traffico in modo efficiente.

Anti-pattern comuni:

- Scelta del tipo di sistema di bilanciatore del carico senza tenere conto dei requisiti del carico di lavoro.
- Mancato utilizzo delle funzionalità del bilanciatore del carico per l'ottimizzazione delle prestazioni.
- Esposizione diretta del carico di lavoro a Internet senza un bilanciatore del carico.
- Instradati tutto il traffico Internet attraverso i bilanciatori del carico esistenti.
- Utilizzi il bilanciamento del TCP carico generico e fai in modo che ogni nodo di calcolo gestisca la crittografia. SSL

Vantaggi dell'adozione di questa best practice: un bilanciatore del carico gestisce il carico variabile del traffico dell'applicazione in una o più zone di disponibilità e consente alta disponibilità, dimensionamento automatico e un migliore utilizzo del carico di lavoro.

Livello di rischio associato se questa best practice non fosse adottata: elevato

Guida all'implementazione

I bilanciatori del carico operano come punto di ingresso per il carico di lavoro, dal quale distribuiscono il traffico alle destinazioni di backend, come istanze di calcolo o container per migliorarne l'utilizzo.

La scelta del tipo corretto di bilanciatore del carico è il primo passaggio per ottimizzare l'architettura. Inizia elencando le caratteristiche del carico di lavoro, ad esempio il protocollo (ad esempio TCPHTTP, TLS, o WebSockets), il tipo di destinazione (come istanze, contenitori o serverless), i requisiti dell'applicazione (come connessioni a lunga durata, autenticazione utente o persistenza) e il posizionamento (come regione, zona locale, Outpost o isolamento zonale).

AWS fornisce diversi modelli per le applicazioni per utilizzare il bilanciamento del carico. [Application Load Balancer](#) è la soluzione ideale per il bilanciamento del carico HTTP e del HTTPS traffico

e fornisce un routing avanzato delle richieste mirato alla distribuzione di architetture applicative moderne, inclusi microservizi e contenitori.

[Network Load Balancer](#) è la soluzione ideale per il bilanciamento del carico del TCP traffico laddove sono richieste prestazioni estreme. È in grado di gestire milioni di richieste al secondo, mantenendo al contempo latenze ridottissime. Inoltre, è ottimizzato per la gestione degli schemi di traffico improvvisi e incostanti.

[Elastic Load Balancing](#) offre la gestione e SSL la TLS decrittografia integrate dei certificati, che ti offrono la flessibilità necessaria per gestire centralmente SSL le impostazioni del sistema di bilanciamento del carico e ridurre il carico di lavoro CPU intensivo.

Dopo aver scelto il bilanciatore del carico appropriato, puoi iniziare a utilizzarne le funzionalità per ridurre la quantità di attività che deve svolgere il backend per distribuire il traffico.

Ad esempio, utilizzando sia Application Load Balancer (ALB) che Network Load Balancer NLB (), è possibile SSL eseguire TLS/encryption offloading, un'opportunità per evitare che CPU l'handshake così TLS impegnativo venga completato dai target e anche per migliorare la gestione dei certificati.

Quando configuri SSL/TLS offloading nel tuo sistema di bilanciamento del carico, quest'ultimo diventa responsabile della crittografia del traffico da e verso i client, mentre consegna il traffico non crittografato ai tuoi backend, liberando le tue risorse di backend e migliorando i tempi di risposta per i client.

Application Load Balancer può anche servire il traffico HTTP /2 senza bisogno di supportarlo sui tuoi obiettivi. Questa semplice decisione può migliorare i tempi di risposta dell'applicazione, poiché HTTP /2 utilizza TCP le connessioni in modo più efficiente.

Nel definire l'architettura, è bene tenere conto dei requisiti di latenza del carico di lavoro. Ad esempio, se un'applicazione è sensibile alla latenza, è possibile scegliere di usare Network Load Balancer, che offre latenze estremamente ridotte. In alternativa, è possibile decidere di avvicinare il carico di lavoro ai clienti sfruttando Application Load Balancer nelle [zone locali AWS](#) o addirittura [AWS Outposts](#).

Un altro aspetto di cui tenere conto per i carichi di lavoro sensibili alla latenza è il bilanciamento del carico tra zone. Con il bilanciamento del carico tra zone, ogni nodo del bilanciatore del carico distribuisce il traffico tra le destinazioni registrate in tutte le zone di disponibilità autorizzate.

Usa Auto Scaling integrato con il bilanciatore del carico. Uno degli aspetti principali di un sistema con prestazioni efficienti riguarda il dimensionamento corretto delle risorse backend. A questo scopo, puoi utilizzare integrazioni dei bilanciatori del carico per le risorse di destinazione backend. Usando l'integrazione dei bilanciatori del carico con gruppi Auto Scaling, le destinazioni vengono aggiunte

o rimosse nel e dal bilanciatore del carico in base alle esigenze, in risposta al traffico in ingresso. I sistemi di bilanciamento del carico possono anche integrarsi con Amazon e [ECSAmazon EKS](#) per carichi di lavoro containerizzati.

- [Amazon ECS - Bilanciamento del carico di servizio](#)
- [Bilanciamento del carico delle applicazioni su Amazon EKS](#)
- [Bilanciamento del carico di rete su Amazon EKS](#)

Passaggi dell'implementazione

- Definisci i tuoi requisiti di bilanciamento del carico, tra cui volume di traffico, disponibilità e scalabilità delle applicazioni.
- Scegli il tipo di sistema di bilanciatore del carico giusto per la tua applicazione.
 - Usa Application Load Balancer per i carichi di lavoro HTTP HTTPS /.
 - Usa Network Load Balancer per HTTP carichi non di lavoro eseguiti su o. TCP UDP
 - Utilizza una combinazione di entrambi ([ALBcome obiettivo di NLB](#)) se desideri sfruttare le funzionalità di entrambi i prodotti. Ad esempio, puoi farlo se desideri utilizzare lo statico IPs di NLB insieme al routing basato sull'HTTPintestazione da ALB o se desideri esporre il tuo HTTP carico di lavoro a un. [AWS PrivateLink](#)
 - [Per un confronto completo dei sistemi di bilanciamento del carico, consulta la sezione Confronto dei prodotti. ELB](#)
- Se possibile, usaSSL/TLSoffloading.
 - ConfiguraHTTPS/TLSlistener con [Application Load Balancer e Network Load Balancer integrati](#) con. [AWS Certificate Manager](#)
 - Tieni presente che alcuni carichi di lavoro potrebbero richiedere la end-to-end crittografia per motivi di conformità. In questo caso, è necessario consentire la crittografia nelle destinazioni.
 - Per le migliori pratiche di sicurezza, consulta [SEC09-BP02](#) Applica la crittografia in transito.
- Seleziona l'algoritmo di routing corretto (solo). ALB
 - L'algoritmo di instradamento può fare la differenza per quanto riguarda l'uso corretto delle destinazioni backend e, di conseguenza, l'impatto sulle prestazioni. Ad esempio, ALB offre [due opzioni per gli algoritmi di routing](#):
 - Numero minimo di richieste in sospenso: usa questa opzione per ottenere una migliore distribuzione del carico nelle destinazioni backend nei casi in cui le richieste per l'applicazione variano per complessità o le destinazioni variano per capacità di elaborazione.

- Round robin: usa questa opzione quando le richieste e le destinazioni sono simili o se devi distribuire equamente le richieste tra le destinazioni.
- Valuta se usare l'isolamento tra zone o quello zonale.
 - Disattiva l'isolamento tra zone (usando l'isolamento zonale) per migliorare la latenza e in caso di domini con errori di zona. È disattivato per impostazione predefinita in NLB e [ALB puoi disattivarlo per gruppo target](#).
 - Attiva l'isolamento tra zone per ottenere disponibilità e flessibilità maggiori. Per impostazione predefinita, l'opzione cross-zone è attivata per ogni gruppo target ALB e [NLB può essere attivata in base al gruppo target](#).
- Attiva HTTP keep-alive per i tuoi HTTP carichi di lavoro (solo). ALB Con questa funzionalità, il load balancer può riutilizzare le connessioni di backend fino alla scadenza del timeout keep-alive, migliorando i tempi di HTTP richiesta e risposta e riducendo anche l'utilizzo delle risorse per gli obiettivi di backend. Per informazioni dettagliate su come eseguire questa operazione per Apache e Nginx, vedi [Quali sono le impostazioni ottimali per l'utilizzo di Apache o come server di backend? NGINX ELB](#)
- Attiva il monitoraggio del tuo bilanciatore del carico.
 - Attiva i log di accesso per [Application Load Balancer](#) e [Network Load Balancer](#).
 - I campi principali da considerare sono, e. ALB request_processing_time response_processing_time
 - I principali campi da prendere in considerazione NLB sono connection_time etls_handshake_time.
 - Preparati a eseguire query sui log quando necessario. [Puoi usare Amazon Athena per interrogare sia i ALB log che i log. NLB](#)
 - [Crea allarmi per metriche relative alle prestazioni, ad esempio per. TargetResponseTime ALB](#)

Risorse

Documenti correlati:

- [ELB confronto tra prodotti](#)
- [AWS Infrastruttura globale](#)
- [Improving Performance and Reducing Cost Using Availability Zone Affinity](#)
- [Step by step for Log Analysis with Amazon Athena](#)
- [Querying Application Load Balancer logs](#)

- [Monitor your Application Load Balancers](#)
- [Monitor your Network Load Balancer](#)
- [Use Elastic Load Balancing to distribute traffic across the instances in your Auto Scaling group](#)

Video correlati:

- [AWS re:Invent 2023: Cosa può fare il networking per la tua applicazione?](#)
- [AWS RE:InForce 20: Come utilizzare Elastic Load Balancing per migliorare il livello di sicurezza su larga scala](#)
- [AWS re:Invent 2018: Elastic Load Balancing: approfondimenti e best practice](#)
- [AWS re:Invent 2021 - Come scegliere il bilanciamento del carico giusto per i tuoi carichi di lavoro AWS](#)
- [AWS re:Invent 2019: ottieni il massimo da Elastic Load Balancing per diversi carichi di lavoro](#)

Esempi correlati:

- [Gateway Load Balancer](#)
- [CDKed AWS CloudFormation esempi per l'analisi dei log con Amazon Athena](#)

PERF04-BP05 Scegli i protocolli di rete per migliorare le prestazioni

Prendi decisioni sui protocolli per la comunicazione tra sistemi e reti in base all'impatto sulle prestazioni del carico di lavoro.

Esiste una relazione tra latenza e larghezza di banda per ottenere il throughput desiderato. Se il trasferimento dei file utilizza Transmission Control Protocol (TCP), latenze più elevate molto probabilmente ridurranno la velocità effettiva complessiva. Esistono approcci per risolvere questo problema con l'ottimizzazione e i protocolli di trasferimento ottimizzati, ma una soluzione consiste nell'utilizzare User Datagram Protocol (UDP).

Anti-pattern comuni:

- Si utilizza TCP per tutti i carichi di lavoro indipendentemente dai requisiti di prestazioni.

Vantaggi dell'adozione di questa best practice: la verifica del protocollo adeguato per la comunicazione tra utenti e componenti del carico di lavoro contribuisce a migliorare l'esperienza

utente complessiva per le applicazioni. Ad esempio, la modalità senza connessione UDP consente l'alta velocità, ma non offre ritrasmissione o alta affidabilità. TCP è un protocollo completo, ma richiede un sovraccarico maggiore per l'elaborazione dei pacchetti.

Livello di rischio associato se questa best practice non fosse adottata: medio

Guida all'implementazione

Se hai la possibilità di scegliere protocolli diversi per la tua applicazione e hai esperienza in questo campo, ottimizza l'applicazione e l'esperienza dell'utente finale utilizzando un protocollo diverso. Tieni conto che questo approccio presenta notevoli difficoltà e dovrebbe essere tentato solo dopo l'ottimizzazione dell'applicazione in altri modi.

Un aspetto principale per il miglioramento delle prestazioni del tuo carico di lavoro consiste nell'identificare i requisiti in termini di latenza e throughput, quindi scegliere i protocolli di rete che ottimizzano le prestazioni.

Quando prendere in considerazione l'utilizzo TCP

TCP fornisce una consegna affidabile dei dati e può essere utilizzato per la comunicazione tra i componenti del carico di lavoro in cui l'affidabilità e la consegna garantita dei dati sono importanti. Molte applicazioni basate sul Web si basano su protocolli TCP basati, come HTTP e HTTPS, per aprire i TCP socket per la comunicazione tra i componenti dell'applicazione. Il trasferimento di dati tramite posta elettronica e file sono applicazioni comuni che utilizzano anch'esse TCP, in quanto si tratta di un meccanismo di trasferimento semplice e affidabile tra i componenti dell'applicazione. L'utilizzo di TLS with TCP può comportare un sovraccarico di comunicazione, con conseguente aumento della latenza e riduzione della velocità effettiva, ma offre anche il vantaggio della sicurezza. Il sovraccarico è dovuto perlopiù al processo di handshake, il cui completamento può richiedere diversi round trip. Al termine del processo di handshake, il sovraccarico dovuto alla crittografia e alla decrittografia dei dati è relativamente ridotto.

Quando prendere in considerazione l'utilizzo UDP

UDP è un connection-less-oriented protocollo ed è quindi adatto per applicazioni che richiedono una trasmissione rapida ed efficiente, come log, monitoraggio e dati VoIP. Inoltre, UDP se disponi di componenti per il carico di lavoro, valuta la possibilità di utilizzare componenti che rispondono a piccole richieste provenienti da un gran numero di client per garantire prestazioni ottimali del carico di lavoro. Datagram Transport Layer Security (DTLS) è l'UDP equivalente di Transport Layer Security (TLS). Quando si utilizza DTLS with UDP, il sovraccarico deriva dalla crittografia e dalla

decriptografia dei dati, poiché il processo di handshake è semplificato. DTLS aggiunge inoltre un piccolo sovraccarico ai UDP pacchetti, in quanto include campi aggiuntivi per indicare i parametri di sicurezza e rilevare eventuali manomissioni.

Quando prendere in considerazione l'utilizzo SRD

Scalable reliable datagram (SRD) è un protocollo di trasporto di rete ottimizzato per carichi di lavoro ad alto throughput grazie alla sua capacità di bilanciare il carico del traffico su più percorsi e ripristinare rapidamente il sistema in caso di interruzioni di pacchetti o errori di collegamento. SRD è quindi utilizzato al meglio per carichi di lavoro di elaborazione ad alte prestazioni (HPC) che richiedono un throughput elevato e una comunicazione a bassa latenza tra i nodi di elaborazione. Possono essere incluse attività di elaborazione in parallelo come la simulazione, la modellazione e l'analisi dei dati che implicano il trasferimento di grandi quantità di dati tra nodi.

Passaggi dell'implementazione

- Utilizzare i servizi [AWS Global Accelerator](#) e [AWS Transfer Family](#) per migliorare il throughput delle applicazioni di trasferimento file online. Il AWS Global Accelerator servizio consente di ridurre la latenza tra i dispositivi client e il carico di lavoro su. AWS Con AWS Transfer Family, puoi utilizzare protocolli TCP basati come Secure Shell File Transfer Protocol (SFTP) e File Transfer Protocol over SSL (FTPS) per scalare e gestire in modo sicuro i trasferimenti di file verso AWS i servizi di archiviazione.
- Utilizza la latenza di rete per determinare se TCP è appropriata per la comunicazione tra i componenti del carico di lavoro. Se la latenza di rete tra l'applicazione client e il server è elevata, l'handshake TCP a tre vie può richiedere del tempo, con un conseguente impatto sulla reattività dell'applicazione. Metriche come time to first byte (TTFB) e round-trip time (RTT) possono essere utilizzate per misurare la latenza di rete. Se il tuo carico di lavoro fornisce contenuti dinamici agli utenti, prendi in considerazione l'utilizzo di [Amazon CloudFront](#), che stabilisce una connessione permanente a ciascuna origine per i contenuti dinamici per rimuovere i tempi di configurazione della connessione che altrimenti rallenterebbero ogni richiesta del client.
- L'utilizzo TLS con TCP o UDP può comportare un aumento della latenza e una riduzione del throughput per il carico di lavoro a causa dell'impatto della crittografia e della decrittografia. Per tali carichi di lavoro, prendi in considerazione l'opzione SSL/TLS offloading su [Elastic Load Balancing](#) per migliorare le prestazioni del carico di lavoro consentendo al sistema di bilanciamento del carico di SSL gestire il processo di crittografia e decrittografia TLS/anziché affidarlo alle istanze di backend. Questo può aiutare a ridurre l'CPU utilizzo delle istanze di backend, migliorando le prestazioni e aumentando la capacità.

- Utilizza [Network Load Balancer \(NLB\)](#) per implementare servizi basati sul UDP protocollo, come autenticazione e autorizzazione, registrazione, DNS IoT e streaming multimediale, per migliorare le prestazioni e l'affidabilità del carico di lavoro. NLB Distribuisce il UDP traffico in entrata su più destinazioni, consentendoti di scalare il carico di lavoro orizzontalmente, aumentare la capacità e ridurre il sovraccarico di un singolo target.
- Per i carichi di lavoro High Performance Computing (HPC), prendi in considerazione l'utilizzo della funzionalità [Elastic Network Adapter \(ENA\) Express](#) che utilizza il SRD protocollo per migliorare le prestazioni di rete fornendo una maggiore larghezza di banda a flusso singolo (25 Gbps) e una latenza di coda inferiore (99,9 percentile) per il traffico di rete tra le istanze. EC2
- Utilizza [Application Load Balancer \(ALB\)](#) per indirizzare e bilanciare il carico del traffico gRPC (Remote Procedure Calls) tra i componenti del carico di lavoro o tra gRPC i client e i servizi. gRPC utilizza il protocollo TCP basato su HTTP /2 per il trasporto e offre vantaggi in termini di prestazioni come un ingombro di rete più leggero, compressione, serializzazione binaria efficiente, supporto per numerose lingue e streaming bidirezionale.

Risorse

Documenti correlati:

- [Come indirizzare il traffico verso Kubernetes UDP](#)
- [Application Load Balancer](#)
- [EC2Rete avanzata su Linux](#)
- [EC2Rete avanzata su Windows](#)
- [EC2Gruppi di collocamento](#)
- [Abilitazione di reti avanzate con Elastic Network Adapter \(ENA\) su istanze Linux](#)
- [Network Load Balancer](#)
- [Prodotti di rete con AWS](#)
- [Transitioning to Latency-Based Routing in Amazon Route 53](#)
- [VPC Endpoint](#)

Video correlati:

- [AWS re:Invent 2022 — Scalabilità delle prestazioni di rete sulle istanze Amazon Elastic Compute Cloud di nuova generazione](#)

- [AWS re:Invent 2022 — Fondamenti del networking delle applicazioni](#)

Esempi correlati:

- [AWS Transit Gateway e soluzioni di sicurezza scalabili](#)
- [Workshop sulle reti AWS](#)

PERF04-BP06 Scegli la posizione del carico di lavoro in base ai requisiti di rete

Valuta le opzioni per il posizionamento delle risorse in modo da diminuire la latenza di rete e migliorare il throughput, fornendo un'esperienza utente ottimale attraverso la riduzione dei tempi di caricamento delle pagine e di trasferimento dei dati.

Anti-pattern comuni:

- Consolidamento di tutte le risorse del carico di lavoro in un'unica posizione geografica.
- Scelta della regione più vicina alla propria posizione, ma non al carico di lavoro dell'utente finale.

Vantaggi dell'adozione di questa best practice: l'esperienza utente è fortemente condizionata dalla latenza tra utente e applicazione. Utilizzando una rete globale appropriata Regioni AWS e AWS privata, è possibile ridurre la latenza e offrire un'esperienza migliore agli utenti remoti.

Livello di rischio associato se questa best practice non fosse adottata: medio

Guida all'implementazione

Le risorse, come le EC2 istanze Amazon, vengono collocate nelle Availability Zones within [Regioni AWS](#), [AWS Local Zones](#) o [AWS Wavelength](#) nelle zone. [AWS Outposts](#) La scelta della posizione influisce su latenza di rete e throughput dall'ubicazione di un utente specifico. I servizi edge come [Amazon CloudFront](#) [AWS Global Accelerator](#) possono essere utilizzati anche per migliorare le prestazioni di rete memorizzando nella cache i contenuti nelle edge location o fornendo agli utenti un percorso ottimale per il carico di lavoro attraverso la rete AWS globale.

Amazon EC2 fornisce gruppi di collocamento per il networking. Un gruppo di collocazione è un raggruppamento logico di istanze per ridurre la latenza. L'utilizzo di gruppi di collocamento con tipi di

istanze supportati e un Elastic Network Adapter (ENA) consente ai carichi di lavoro di partecipare a una rete a 25 Gbps a bassa latenza e con jitter ridotto. I gruppi di collocazione sono consigliati per i carichi di lavoro che traggono beneficio da reti a bassa latenza, throughput di rete elevato o entrambi.

[I servizi sensibili alla latenza vengono forniti nelle sedi periferiche utilizzando una rete AWS globale, come Amazon CloudFront](#) Queste edge location forniscono in genere servizi come Content Delivery Network (CDN) e Domain Name System (). DNS Disponendo di questi servizi all'edge, i carichi di lavoro possono rispondere con bassa latenza alle richieste di contenuto o DNS risoluzione. Inoltre, possono offrire servizi geografici come la geotargetizzazione dei contenuti (ossia fornire contenuti diversi in base alla posizione dell'utente finale) o l'instradamento basato sulla latenza, per indirizzare gli utenti alla regione più vicina (latenza minima).

Usa i servizi edge per ridurre la latenza e abilitare la memorizzazione nella cache dei contenuti. Configura correttamente il controllo della cache per entrambi DNS eHTTP/HTTPSper ottenere il massimo vantaggio da questi approcci.

Passaggi dell'implementazione

- Acquisisci informazioni sul traffico IP in entrata e in uscita dalle interfacce di rete.
 - [Registrazione del traffico IP utilizzando VPC Flow Logs](#)
 - [Come viene preservato l'indirizzo IP del client in AWS Global Accelerator](#)
- Analizza i modelli di accesso alla rete nel tuo carico di lavoro per capire come gli utenti usano la tua applicazione.
 - Utilizza strumenti di monitoraggio, come [Amazon CloudWatch](#) e [AWS CloudTrail](#), per raccogliere dati sulle attività di rete.
 - Analizza i dati per identificare il modello di accesso alla rete.
- Seleziona regioni appropriate per l'implementazione del carico di lavoro in base ai seguenti elementi chiave:
 - Ubicazione dei dati per le applicazioni a uso intensivo di dati, ad esempio applicazioni di big data e machine learning, il codice dell'applicazione dovrebbe essere eseguito il più vicino possibile ai dati.
 - Ubicazione degli utenti: per le applicazioni rivolte agli utenti, scegli una regione o più regioni vicine agli utenti del carico di lavoro.
 - Altri vincoli: prendi in considerazione vincoli come costi e conformità, come illustrato in [What to Consider when Selecting a Region for your Workloads.](#)

- Usa le [zone locali AWS](#) per eseguire carichi di lavoro come il rendering video. Le zone locali consentono di sfruttare i vantaggi derivanti dalla disponibilità di risorse di calcolo e archiviazione più vicine agli utenti finali.
- Usa [AWS Outposts](#) per carichi di lavoro che devono rimanere on-premises, ma vuoi che vengano eseguiti in modo ottimale con il resto degli altri carichi di lavoro in AWS.
- Applicazioni come lo streaming video in diretta ad alta risoluzione, l'audio ad alta fedeltà e la realtà aumentata o la realtà virtuale (AR/VR) richiedono dispositivi 5G. ultra-low-latency Per tali applicazioni, considera. [AWS Wavelength](#) AWS Wavelength incorpora servizi di AWS elaborazione e archiviazione nelle reti 5G, fornendo un'infrastruttura di edge computing mobile per lo sviluppo, l'implementazione e la scalabilità delle applicazioni. ultra-low-latency
- Usa la cache locale o le [soluzioni di caching AWS](#) per i dati di frequente utilizzo per migliorare le performance, ridurre lo spostamento dei dati e minimizzare l'impatto ambientale.

Servizio	Quando usare
Amazon CloudFront	Utilizzalo per memorizzare nella cache contenuti statici come immagini, script e video, nonché contenuti dinamici come API risposte o applicazioni web.
Amazon ElastiCache	Usalo per memorizzare nella cache i contenuti per le applicazioni Web.
DynamoDB Accelerator	Usalo per aggiungere accelerazione in memoria alle tabelle DynamoDB.

- Utilizza servizi in grado di supportarti nell'esecuzione del codice in posizioni più vicine agli utenti del carico di lavoro, come i seguenti:

Servizio	Quando usare
Lambda@Edge	Usalo per operazioni a uso intensivo di risorse di calcolo eseguite quando gli oggetti non si trovano nella cache.
CloudFront Funzioni Amazon	Utilizzalo per casi d'uso semplici come richieste HTTP (s) o manipolazioni di risposte

Servizio	Quando usare
AWS IoT Greengrass	Usale per eseguire la memorizzazione nella cache di risorse di calcolo, messaggistica e dati per i dispositivi connessi.

- Alcune applicazioni richiedono punti di ingresso fissi o prestazioni più elevate attraverso la riduzione della latenza di ricezione del primo byte e l'instabilità e l'aumento del throughput. Queste applicazioni possono trarre vantaggio da servizi di rete che forniscono indirizzi IP anycast statici e TCP terminazioni in postazioni periferiche. [AWS Global Accelerator](#) possono migliorare le prestazioni delle applicazioni fino al 60% e fornire un failover rapido per architetture multiregionali. AWS Global Accelerator fornisce indirizzi IP anycast statici che fungono da punto di ingresso fisso per le applicazioni ospitate in una o più applicazioni. Regioni AWS Questi indirizzi IP consentono al traffico di entrare nella rete AWS globale il più vicino possibile agli utenti. AWS Global Accelerator riduce il tempo di configurazione iniziale della connessione stabilendo una TCP connessione tra il client e la AWS edge location più vicina al client. Rivedi l'utilizzo di AWS Global Accelerator per migliorare le prestazioni dei tuoi UDP carichi di lavoro TCP/e fornire un failover rapido per architetture multiregionali.

Risorse

Best practice correlate:

- [COST07-BP02 Implementazione delle regioni in base ai costi](#)
- [COST08-BP03 Implementazione di servizi per ridurre i costi di trasferimento dei dati](#)
- [REL10-BP01 Implementa il carico di lavoro in più sedi](#)
- [REL10-BP02 Seleziona le posizioni appropriate per l'implementazione in più sedi](#)
- [SUS01-BP01 Scegli la regione in base ai requisiti aziendali e agli obiettivi di sostenibilità](#)
- [SUS02-BP04 Ottimizza il posizionamento geografico dei carichi di lavoro in base ai requisiti di rete](#)
- [SUS04-BP07 Riduci al minimo lo spostamento dei dati tra le reti](#)

Documenti correlati:

- [AWS Infrastruttura globale](#)
- [AWS Local Zones e AWS Outposts scelta della tecnologia giusta per il tuo carico di lavoro edge](#)
- [Placement groups](#)
- [AWS Local Zones](#)
- [AWS Outposts](#)
- [AWS Wavelength](#)
- [Amazon CloudFront](#)
- [AWS Global Accelerator](#)
- [AWS Direct Connect](#)
- [AWS Site-to-Site VPN](#)
- [Amazon Route 53](#)

Video correlati:

- [AWS Video esplicativo su Local Zones](#)
- [AWS Outposts: Overview and How it Works](#)
- [AWS re:Invent 2023 - Una strategia di migrazione per carichi di lavoro edge e locali](#)
- [AWS re:Invent 2021 -: Portare l'esperienza in sede AWS OutpostsAWS](#)
- [AWS re:Invent 2020:: Esegui app con latenza AWS Wavelength ultra bassa su 5G Edge](#)
- [AWS re:Invent 2022 - AWS Local Zones: creazione di applicazioni per un edge distribuito](#)
- [AWS re:Invent 2021 - Creazione di siti Web a bassa latenza con Amazon CloudFront](#)
- [AWS re:Invent 2022 - Migliora le prestazioni e la disponibilità con AWS Global Accelerator](#)
- [AWS re:Invent 2022 - Costruisci la tua rete WAN utilizzando AWS](#)
- [AWS re:Invent 2020: gestione globale del traffico con Amazon Route 53](#)

Esempi correlati:

- [AWS Global Accelerator Workshop sul routing personalizzato](#)
- [Handling Rewrites and Redirects using Edge Functions](#)

PERF04-BP07 Ottimizza la configurazione di rete in base a metriche

Usa i dati raccolti e analizzati per prendere decisioni informate riguardo l'ottimizzazione della configurazione della tua rete.

Anti-pattern comuni:

- Si ritiene che tutti i problemi relativi alle prestazioni siano correlati all'applicazione.
- Verifica delle prestazioni di rete solo da una posizione vicina a quella in cui è stato distribuito il carico di lavoro.
- Uso di configurazioni predefinite per tutti i servizi di rete.
- Provisioning in eccesso di risorse di rete per fornire capacità sufficiente.

Vantaggi dell'adozione di questa best practice: la raccolta delle metriche necessarie per la rete AWS e l'implementazione di strumenti di monitoraggio di rete permettono di identificare le prestazioni di rete e ottimizzare le configurazioni di rete.

Livello di rischio associato se questa best practice non fosse adottata: basso

Guida all'implementazione

Il monitoraggio del traffico da e verso VPCs, le sottoreti o le interfacce di rete è fondamentale per capire come utilizzare le risorse di rete e ottimizzare le configurazioni di rete. AWS Utilizzando i seguenti strumenti AWS di rete, è possibile esaminare ulteriormente le informazioni sull'utilizzo del traffico, sull'accesso alla rete e sui registri.

Passaggi dell'implementazione

- Identifica le metriche prestazionali chiave, come la latenza o la perdita di pacchetti, da raccogliere. AWS fornisce diversi strumenti che possono aiutarti a raccogliere queste metriche. Usando i seguenti strumenti, puoi esaminare ulteriormente le informazioni sull'utilizzo del traffico, sull'accesso alla rete e sui log:

AWS strumento	Dove usarlo
Gestione indirizzi VPC IP di Amazon.	IPAM Utilizzalo per pianificare, tracciare e monitorare gli indirizzi IP per i tuoi carichi di lavoro AWS e quelli locali. Si tratta di una best practice per ottimizzare l'utilizzo e l'allocatione degli indirizzi IP.
VPC Registri di flusso	Usa VPC Flow Logs per acquisire informazioni dettagliate sul traffico da e verso le interfacce di rete del tuo VPC. Con VPC Flow Logs, puoi diagnosticare regole di gruppo di sicurezza eccessivamente restrittive o permissive e determinare la direzione del traffico da e verso le interfacce di rete.
AWS Transit Gateway Log di flusso	Utilizzate AWS Transit Gateway Flow Logs per acquisire informazioni sul traffico IP in entrata e in uscita dai gateway di transito.
DNS registrazione delle interrogazioni	Registra le informazioni sulle DNS interrogazioni pubbliche o private ricevute da Route 53. Con DNS i log, è possibile ottimizzare le DNS configurazioni comprendendo il dominio o il sottodominio richiesto o le EDGE posizioni di Route 53 che hanno risposto alle query. DNS
Reachability Analyzer	Con Reachability Analyzer puoi analizzare la raggiungibilità della rete ed eseguirne il debug. Reachability Analyzer è uno strumento di analisi della configurazione che consente di eseguire test di connettività tra una risorsa di origine e una risorsa di destinazione nel proprio VPC. Lo strumento in questione consente di verificare la corrispondenza fra configurazione e connettività desiderata.

AWS strumento	Dove usarlo
Strumento di analisi degli accessi alla rete	È possibile utilizzare lo Strumento di analisi degli accessi alla rete per comprendere l'accesso di rete alle risorse. Puoi usare lo Strumento di analisi degli accessi alla rete per specificare i requisiti di accesso alla rete e identificare i potenziali percorsi di rete che non li soddisfano. Ottimizzando la configurazione di rete corrispondente, puoi determinare e verificare lo stato della rete e indicare se la rete su AWS soddisfa i requisiti di conformità.
Amazon CloudWatch	Usa Amazon CloudWatch e attiva i parametri appropriati per le opzioni di rete. Assicurati di scegliere le metriche di rete corrette per il carico di lavoro. Ad esempio, puoi attivare le metriche per VPC Network Address Usage, VPC NAT Gateway AWS Transit Gateway, VPN tunnel AWS Network Firewall, Elastic Load Balancing e. AWS Direct Connect Il monitoraggio continuo delle metriche è una procedura utile per osservare e identificare lo stato e l'utilizzo della rete che semplifica l'ottimizzazione della configurazione di rete in base alle osservazioni.

AWS strumento	Dove usarlo
AWS Network Manager	Utilizzando AWS Network Manager, è possibile monitorare le prestazioni in tempo reale e cronologiche della rete AWS globale per scopi operativi e di pianificazione. Network Manager fornisce la latenza di rete aggregata tra le zone di disponibilità Regioni AWS e all'interno di ciascuna zona di disponibilità, consentendovi di comprendere meglio in che modo le prestazioni delle applicazioni si relazionano con le prestazioni della rete sottostante. AWS
Amazon CloudWatch RUM	Usa Amazon CloudWatch RUM per raccogliere le metriche che ti forniscono le informazioni che ti aiutano a identificare, comprendere e migliorare l'esperienza utente.

- Identifica i principali oratori e i modelli di traffico delle applicazioni utilizzando VPC e AWS Transit Gateway Flow Logs.
- Valuta e ottimizza la tua attuale architettura di rete VPCs, inclusi sottoreti e routing. Ad esempio, potete valutare la diversità del VPC peering o AWS Transit Gateway aiutarvi a migliorare il networking nella vostra architettura.
- Valuta i percorsi di instradamento nella tua rete per verificare che venga sempre utilizzato il percorso più breve tra le destinazioni. Lo Strumento di analisi degli accessi alla rete è utile in questa operazione.

Risorse

Documenti correlati:

- [DNS Registrazione pubblica delle interrogazioni](#)
- [Che cos'è IPAM?](#)
- [What is Reachability Analyzer?](#)
- [What is Network Access Analyzer?](#)

- [CloudWatchmetriche per te VPCs](#)
- [Ottimizza le prestazioni e riduci i costi per l'analisi di rete con VPC Flow Logs in formato Apache Parquet](#)
- [Monitoraggio delle tue reti globali e principali con i CloudWatch parametri di Amazon](#)
- [Continuously monitor network traffic and resources](#)

Video correlati:

- [AWS re:Invent 2023: guida per sviluppatori al cloud networking](#)
- [AWS re:Invent 2023 — Pronti per il futuro? Designing networks for growth and flexibility](#)
- [AWS re:Invent 2023 — Design avanzati e nuove funzionalità VPC](#)
- [AWS re:Invent 2022 — Approfondisci l'infrastruttura di rete AWS](#)
- [AWS re:Invent 2020 — Le migliori pratiche e suggerimenti per il networking con il Well-Architected Framework AWS](#)
- [AWS re:Invent 2020 — Monitoraggio e risoluzione dei problemi del traffico di rete](#)

Esempi correlati:

- [Workshop sulle reti AWS](#)
- [AWS Network Monitoring](#)
- [Osservazione e diagnosi della rete su AWS](#)
- [Individuazione e risoluzione degli errori di configurazione della rete su AWS](#)

Processo e cultura

Durante la fase di progettazione dei carichi di lavoro, esistono principi e pratiche che è possibile adottare per gestire al meglio carichi di lavoro cloud efficienti e ad alte prestazioni. Questa area di interesse offre le best practice per aiutarti ad adottare una cultura che promuova l'efficienza delle prestazioni dei carichi di lavoro cloud.

Per sviluppare questa cultura, considera questi principi chiave:

- **Infrastructure as code:** definisci il tuo modello Infrastructure as code tramite approcci come i modelli di AWS CloudFormation. L'uso dei modelli ti consente di collocare la tua infrastruttura nel controllo sorgente, insieme al codice e alle configurazioni dell'applicazione. Ciò ti permette di applicare le stesse procedure di sviluppo software all'infrastruttura, in modo da accelerare l'iterazione.
- **Pipeline di implementazione:** usa una pipeline di integrazione continua/implementazione continua (CI/CD) (ad esempio, repository del codice sorgente, sistemi di sviluppo, distribuzione e automazione dei test) per distribuire la tua infrastruttura. Ciò ti consente di effettuare l'implementazione in modo ripetibile, coerente ed economicamente vantaggioso nel corso dell'iterazione.
- **Parametri ben definiti:** configura e monitora le metriche per raccogliere gli indicatori chiave di prestazione (KPI). Ti consigliamo di adottare parametri tecnici e aziendali. Per i siti Web o le app mobili, le metriche principali sono il tempo di acquisizione al primo byte o il rendering. Gli altri parametri generalmente validi includono il numero di thread, il tasso di rimozione di oggetti inutili (garbage collection) e gli stati di attesa. I parametri aziendali, come il costo cumulativo aggregato per richiesta, possono indicarti due modi per ridurre i costi. Valuta attentamente il modo in cui prevedi di interpretare i parametri. Ad esempio, potresti scegliere il 99° percentile o quello massimo anziché il valore medio.
- **Automatizza i test delle prestazioni:** nell'ambito del processo di implementazione, avvia automaticamente i test delle prestazioni dopo che quelli dall'esecuzione più rapida hanno dato esito positivo. L'automazione deve creare un nuovo ambiente, configurare le condizioni iniziali come i dati del test ed eseguire una serie di benchmark e test di carico. I risultati dei test devono essere confrontati con la build, in modo da monitorare le variazioni delle prestazioni nel corso del tempo. Per i test di lunga durata, puoi inserirli nella pipeline in maniera asincrona rispetto al resto della build. In alternativa, puoi eseguire i test delle prestazioni negli orari notturni, tramite le istanze spot di Amazon EC2.
- **Generazione del carico:** crea una serie di script di test che replichino i percorsi utente sintetici o pre-registrati. Tali script devono essere idempotenti e non devono essere associati in coppie.

Inoltre, potrebbe essere necessario includere script preliminari per garantire risultati validi. Testa gli script il più possibile, per assicurarti che replichino le abitudini di utilizzo in produzione. Puoi usare soluzioni software as a service (SaaS) per generare il carico. Valuta se l'utilizzo delle soluzioni [Marketplace AWS](#) e le [istanze spot](#) possono essere modi convenienti per generare il carico.

- **Visibilità delle prestazioni:** i parametri principali devono essere visibili dal team, in particolar modo quelli relativi a ciascuna versione della build. Ciò ti consente di rilevare tendenze positive o negative rilevanti nel corso del tempo. Dovresti anche visualizzare i parametri sul numero di errori o eccezioni per assicurarti di testare un sistema funzionante.
- **Visualizzazione:** sfrutta le tecniche di visualizzazione che indicano in modo chiaro i punti in cui si verificano problemi di prestazioni, hot spot, stati di attesa o utilizzo ridotto. Sovrapponi i parametri delle prestazioni ai diagrammi architetturali: i grafici delle chiamate o il codice possono aiutarti a individuare più rapidamente i problemi.
- **Revisione regolare dei processi:** le prestazioni scarse delle architetture sono in genere il risultato di un processo di revisione delle prestazioni inesistente o incompleto. Se la tua architettura offre prestazioni insufficienti, l'implementazione di un processo di revisione delle prestazioni ti consente di favorire il miglioramento delle iterazioni.
- **Ottimizzazione continua:** adotta una cultura per ottimizzare continuamente l'efficienza delle prestazioni del tuo carico di lavoro cloud.

Best practice

- [PERF05-BP01 Stabilire indicatori chiave di prestazione \(KPIs\) per misurare lo stato e le prestazioni del carico di lavoro](#)
- [PERF05-BP02 Utilizza soluzioni di monitoraggio per comprendere le aree in cui le prestazioni sono più critiche](#)
- [PERF05-BP03 Definire un processo per migliorare le prestazioni del carico di lavoro](#)
- [PERF05-BP04 Load Esegui un test del tuo carico di lavoro](#)
- [PERF05-BP05 Usa l'automazione per risolvere in modo proattivo i problemi relativi alle prestazioni](#)
- [PERF05-BP06 Conserva il carico di lavoro e i servizi up-to-date](#)
- [PERF05-BP07 Rivedi le metriche a intervalli regolari](#)

PERF05-BP01 Stabilire indicatori chiave di prestazione (KPIs) per misurare lo stato e le prestazioni del carico di lavoro

Identifica quelli KPIs che misurano quantitativamente e qualitativamente le prestazioni del carico di lavoro. KPIsti aiutano a misurare lo stato e le prestazioni di un carico di lavoro correlato a un obiettivo aziendale.

Anti-pattern comuni:

- Monitoraggio dei parametri a livello di sistema solo per avere una visione del carico di lavoro e mancata valutazione degli impatti aziendali di tali parametri.
- Partite dal presupposto che i vostri dati KPIs siano già stati pubblicati e condivisi come dati metrici standard.
- Non definisci un valore quantitativo e misurabile. KPI
- Non siete in linea KPIs con gli obiettivi o le strategie aziendali.

Vantaggi derivanti dall'adozione di questa best practice: l'identificazione di specifiche KPIs che rappresentino lo stato e le prestazioni del carico di lavoro aiuta ad allineare i team sulle loro priorità e a definire risultati aziendali di successo. La condivisione di tali metriche con tutti i reparti fornisce visibilità e allineamento su soglie, aspettative e impatto aziendale.

Livello di rischio associato se questa best practice non fosse adottata: elevato

Guida all'implementazione

KPIsconsentire ai team aziendali e tecnici di allinearsi sulla misurazione degli obiettivi e delle strategie e sul modo in cui questi fattori si combinano per produrre risultati aziendali. Ad esempio, il carico di lavoro di un sito Web può utilizzare il tempo di caricamento della pagina come indicazione delle prestazioni complessive. Questa metrica sarebbe uno dei vari punti dati che misurano l'esperienza dell'utente. Oltre a identificare le soglie di tempo di caricamento della pagina, occorre documentare il risultato atteso o il rischio aziendale in caso di mancato raggiungimento delle prestazioni ideali. Un lungo tempo di caricamento della pagina si ripercuote direttamente sugli utenti finali, peggiora la loro esperienza d'uso e può portare a una perdita di clienti. Quando definisci le tue KPI soglie, combina i benchmark di settore e le aspettative degli utenti finali. Ad esempio, se l'attuale benchmark di settore è una pagina Web che viene caricata entro un periodo di tempo di due secondi, ma gli utenti finali si aspettano che una pagina Web venga caricata entro un periodo di tempo di un secondo, è necessario prendere in considerazione entrambi questi dati quando si stabilisce il KPI

Il team deve valutare il carico di lavoro KPIs utilizzando dati granulari e storici in tempo reale come riferimento e creare dashboard che eseguano calcoli metrici sui dati per ricavare informazioni operative e di utilizzo. KPIs devono essere documentati e includere soglie che supportano gli obiettivi e le strategie aziendali e devono essere mappati alle metriche monitorate. KPIs devono essere riesaminate quando gli obiettivi aziendali, le strategie o i requisiti degli utenti finali cambiano.

Passaggi dell'implementazione

- **Identifica le parti interessate:** identifica e documenta le principali parti interessate aziendali, compresi i team di sviluppo e operativi.
- **Definisci gli obiettivi:** collabora con queste parti interessate per definire e documentare gli obiettivi del carico di lavoro. Considera gli aspetti critici relativi alle prestazioni dei carichi di lavoro, come il throughput, i tempi di risposta e i costi, nonché gli obiettivi aziendali, come la soddisfazione degli utenti.
- **Esamina le best practice del settore:** esamina le best practice del settore per identificare le soluzioni pertinenti in KPIs linea con gli obiettivi del carico di lavoro.
- **Individua le metriche:** identifica le metriche in linea con gli obiettivi del carico di lavoro e in grado di aiutarti a misurare prestazioni e obiettivi aziendali. Stabilisci KPIs in base a queste metriche. ad esempio le misurazioni del tempo medio di risposta o del numero di utenti simultanei.
- **Definisci e documenta KPIs:** utilizza le migliori pratiche del settore e gli obiettivi del carico di lavoro per fissare obiettivi per il tuo carico di lavoro. KPI Utilizza queste informazioni per impostare KPI soglie di gravità o livello di allarme. Identifica e documenta il rischio e l'impatto di un problema KPI non soddisfatto.
- **Implementa il monitoraggio:** utilizza strumenti di monitoraggio come [Amazon CloudWatch](#) o [AWS Config](#) per raccogliere metriche e misurare KPIs.
- **Comunicazione visiva KPIs:** utilizza strumenti di dashboard come [Amazon QuickSight](#) per visualizzare e comunicare KPIs con le parti interessate.
- **Analizza e ottimizza:** rivedi e analizza regolarmente KPIs per identificare le aree del tuo carico di lavoro che devono essere migliorate. Collabora con le parti interessate per implementare tali miglioramenti.
- **Rivedi e perfeziona:** rivedi regolarmente le metriche e KPIs valutane l'efficacia, soprattutto quando gli obiettivi aziendali o le prestazioni del carico di lavoro cambiano.

Risorse

Documenti correlati:

- [CloudWatchdocumentazione](#)
- [Monitoraggio, registrazione e prestazioni s AWS Partner](#)
- [AWS strumenti di osservabilità](#)
- [L'importanza degli indicatori chiave di prestazione \(KPIs\) per le migrazioni al cloud su larga scala](#)
- [Come monitorare l'ottimizzazione dei costi KPIs con la Dashboard KPI](#)
- [Documentazione di X-Ray](#)
- [Utilizzo delle CloudWatch dashboard di Amazon](#)
- [Amazon QuickSight KPIs](#)

Video correlati:

- [AWS re:Invent 2023 - Ottimizza costi e prestazioni e monitora i progressi verso la mitigazione](#)
- [AWS re:Invent 2023 - Gestisci gli eventi del ciclo di vita delle risorse su larga scala con AWS Health](#)
- [AWS re:Invent 2023 - Prestazioni ed efficienza su Pinterest: ottimizzazione delle istanze più recenti](#)
- [AWS re:Invent 2022 - ottimizzazione: misure attuabili per risultati immediati AWS](#)
- [AWS re:Invent 2023 - Costruire un'efficace strategia di osservabilità](#)
- [AWS Summit SF 2022 - Osservabilità completa e monitoraggio delle applicazioni con AWS](#)
- [AWS re:Invent 2023 - Scalabilità per i primi 10 milioni di utenti AWS](#)
- [AWS re:Invent 2022 - In che modo Amazon utilizza metriche migliori per migliorare le prestazioni del sito Web](#)
- [Creazione di una strategia di metrica efficace per la tua azienda | Eventi AWS](#)

Esempi correlati:

- [Creazione di una dashboard con Amazon QuickSight](#)

PERF05-BP02 Utilizza soluzioni di monitoraggio per comprendere le aree in cui le prestazioni sono più critiche

Comprendi e identifica le aree in cui l'aumento delle prestazioni del carico di lavoro determinerà un impatto positivo sull'efficienza o sull'esperienza del cliente. Ad esempio, un sito web che ha una grande quantità di interazione con i clienti può trarre vantaggio dall'utilizzo dei servizi edge per spostare la distribuzione di contenuti più vicino ai clienti.

Anti-pattern comuni:

- Si presume che le metriche di elaborazione standard come l'CPUutilizzo o la pressione della memoria siano sufficienti per individuare problemi di prestazioni.
- Utilizzo solo dei parametri predefiniti registrati dal software di monitoraggio selezionato.
- Revisione dei parametri solo quando c'è un problema.

Vantaggi derivanti dall'adozione di questa best practice: la comprensione delle aree critiche delle prestazioni aiuta i proprietari dei carichi di lavoro a monitorare KPIs e dare priorità ai miglioramenti ad alto impatto.

Livello di rischio associato se questa best practice non fosse adottata: elevato

Guida all'implementazione

Imposta il end-to-end tracciamento per identificare i modelli di traffico, la latenza e le aree critiche di prestazioni. Monitora gli schemi di accesso ai dati per query lente o dati scarsamente frammentati e partizionati. Identifica le aree vincolate del carico di lavoro utilizzando test o monitoraggio del carico.

Aumenta l'efficienza delle prestazioni esaminando l'architettura, gli schemi di traffico e gli schemi di accesso ai dati e identifica la latenza e i tempi di elaborazione. Identifica i potenziali colli di bottiglia che potrebbero influire sull'esperienza del cliente man mano che il carico di lavoro aumenta. Dopo aver identificato queste aree, individua quale soluzione puoi implementare per evitare tali problemi di prestazioni.

Passaggi dell'implementazione

- Imposta il end-to-end monitoraggio per acquisire tutti i componenti e le metriche del carico di lavoro. Ecco alcuni esempi di soluzioni di monitoraggio su AWS

Servizio	Dove usarlo
Amazon CloudWatch Real-User Monitoring (RUM)	Per acquisire i parametri delle prestazioni delle applicazioni da sessioni lato client e frontend di utenti reali.
AWS X-Ray	Per tenere traccia del traffico nei livelli dell'applicazione e identificare la latenza tra componenti e dipendenze. Utilizza le mappe del servizio X-Ray per osservare le relazioni e la latenza tra i componenti del carico di lavoro.
Informazioni dettagliate sulle prestazioni del servizio Amazon Relational Database	Per osservare i parametri delle prestazioni del database e identificare le prestazioni da migliorare.
Monitoraggio RDS avanzato di Amazon	Per osservare i parametri delle prestazioni del sistema operativo del database.
Amazon DevOps Guru	Per rilevare modelli operativi anomali in modo da poter identificare i problemi operativi prima che abbiano un impatto sui clienti.

- Esegui i test per generare parametri, identificare schemi di traffico, colli di bottiglia e aree con prestazioni critiche. Ecco alcuni esempi di come eseguire i test:
 - Configura [CloudWatchSynthetic Canaries](#) per imitare le attività degli utenti basate su browser in modo programmatico utilizzando cron job Linux o espressioni di frequenza per generare metriche coerenti nel tempo.
 - Usa la soluzione [Test di carico distribuito di AWS](#) per generare picchi di traffico o testare il carico di lavoro al tasso di crescita previsto.
- Valuta parametri e dati di telemetria per identificare le aree critiche delle prestazioni. Esamina queste aree con il tuo team per determinare il monitoraggio e le soluzioni per evitare i colli di bottiglia.
- Sperimenta i miglioramenti delle prestazioni e valuta tali modifiche con i dati. Ad esempio, puoi usare [CloudWatchEvidently](#) per testare nuovi miglioramenti e impatti prestazionali sul tuo carico di lavoro.

Risorse

Documenti correlati:

- [Cosa c'è di nuovo in AWS Observability a re:Invent 2023](#)
- [Amazon Builders' Library](#)
- [Documentazione di X-Ray](#)
- [Amazon CloudWatch RUM](#)
- [Amazon DevOps Guru](#)

Video correlati:

- [AWS re:Invent 2023 - \[LAUNCH\] Monitoraggio delle applicazioni per carichi di lavoro moderni](#)
- [AWS re:Invent 2023 - Implementazione dell'osservabilità delle applicazioni](#)
- [AWS re:Invent 2023 - Costruire una strategia di osservabilità efficace](#)
- [AWS Summit SF 2022 - Osservabilità completa e monitoraggio delle applicazioni con AWS](#)
- [AWS Re:Invent 2022 - AWS ottimizzazione: passaggi attuabili per risultati immediati](#)
- [AWS re:Invent 2022 - The Amazon Builders' Library: 25 anni di eccellenza operativa di Amazon](#)
- [AWS re:Invent 2022 - In che modo Amazon utilizza metriche migliori per migliorare le prestazioni del sito Web](#)
- [Monitoraggio visivo delle applicazioni con Amazon CloudWatch Synthetics](#)

Esempi correlati:

- [Misura il tempo di caricamento della pagina con Amazon CloudWatch Synthetics](#)
- [Client CloudWatch RUM Web Amazon](#)
- [X-Ray SDK per Python](#)
- [Test di carico distribuito su AWS](#)

PERF05-BP03 Definire un processo per migliorare le prestazioni del carico di lavoro

Definisci un processo per valutare i nuovi servizi, i modelli di progettazione, i tipi di risorse e le configurazioni man mano che diventano disponibili. Ad esempio, esegui test delle prestazioni esistenti sulle nuove offerte di istanze per determinare il loro potenziale per migliorare il carico di lavoro.

Anti-pattern comuni:

- Si ritiene che l'architettura corrente diventi statica e non venga aggiornata nel corso del tempo.
- Introduzione di modifiche all'architettura nel tempo senza dei parametri che le giustifichino.

Vantaggi dell'adozione di questa best practice: definire un processo per apportare modifiche all'architettura consente ai dati raccolti di influenzare la progettazione del carico di lavoro nel corso del tempo.

Livello di rischio associato se questa best practice non fosse adottata: medio

Guida all'implementazione

Le prestazioni del carico di lavoro presentano alcuni vincoli principali. Documentali, in modo da sapere quali tipi di innovazione potrebbero migliorare le prestazioni del carico di lavoro. Utilizza queste informazioni quando vieni a conoscenza di nuovi servizi o tecnologie, man mano che si rendono disponibili, in modo da identificare le soluzioni per ovviare ai vincoli o ai colli di bottiglia.

Determina i principali vincoli riguardanti le prestazioni del carico di lavoro. Documenta i vincoli prestazionali del carico di lavoro in modo da sapere quali tipi di innovazione potrebbero migliorare le prestazioni del carico di lavoro.

Passaggi dell'implementazione

- Identificazione KPIs: identifica le prestazioni del carico di lavoro KPIs come indicato nella tabella di base del carico di lavoro. [PERF05-BP01 Stabilire indicatori chiave di prestazione \(KPIs\) per misurare lo stato e le prestazioni del carico di lavoro](#)
- Implementa il monitoraggio: utilizza [strumenti di AWS osservabilità](#) per raccogliere metriche e misurare le prestazioni. KPIs
- Effettua analisi: conduci analisi approfondite per individuare le aree (come la configurazione e il codice applicativo) del carico di lavoro con prestazioni insufficienti, come indicato in [PERF05-](#)

[BP02 Utilizza soluzioni di monitoraggio per comprendere le aree in cui le prestazioni sono più critiche](#). Usa i tuoi strumenti di analisi e prestazioni per individuare la strategia di miglioramento delle prestazioni.

- Convalida i miglioramenti: utilizza gli ambienti sandbox o di preproduzione per convalidare l'efficacia della strategia di miglioramento.
- Implementa le modifiche: implementa le modifiche nella produzione e monitora in modo continuo le prestazioni del carico di lavoro. Documenta i miglioramenti e comunica i risultati alle parti interessate.
- Riesamina e perfeziona: rivedi con regolarità il processo di miglioramento delle prestazioni per individuare le aree di miglioramento.

Risorse

Documenti correlati:

- [Blog AWS](#)
- [Cosa c'è di nuovo con AWS](#)
- [AWS Skill Builder](#)

Video correlati:

- [AWS re:Invent 2022 - Fornire architetture sostenibili e ad alte prestazioni](#)
- [AWS re:Invent 2023 - Ottimizza costi e prestazioni e monitora i progressi verso la mitigazione](#)
- [AWS re:Invent 2022 - AWS ottimizzazione: misure attuabili per risultati immediati](#)
- [AWS re:Invent 2022 - Ottimizza i tuoi carichi di lavoro seguendo le migliori pratiche AWS](#)

Esempi correlati:

- [AWS Github](#)

PERF05-BP04 Load Esegui un test del tuo carico di lavoro

Esegui il test del carico di lavoro per verificare che sia in grado di gestire il carico di produzione e individuare eventuali colli di bottiglia nelle prestazioni.

Anti-pattern comuni:

- Test delle singole parti del carico di lavoro, ma non dell'intero carico di lavoro.
- Test di carico eseguito su un'infrastruttura diversa dall'ambiente di produzione.
- Test di carico eseguiti solo per il carico previsto e non oltre, per prevedere dove si potrebbero riscontrare problemi futuri.
- Esegui test di carico senza consultare la [Amazon EC2 Testing Policy](#) e inviare un modulo di invio di eventi simulati. Ciò comporta la mancata esecuzione del test, in quanto sembra un evento. denial-of-service

Vantaggi dell'adozione di questa best practice: misurando le prestazioni in un test di carico, potrai vedere dove avrà luogo l'impatto con l'aumento del carico. In questo modo puoi anticipare le modifiche necessarie prima che influiscano sul carico di lavoro.

Livello di rischio associato se questa best practice non fosse adottata: basso

Guida all'implementazione

Il test di carico nel cloud è un processo volto a misurare le prestazioni del carico di lavoro in condizioni realistiche e con il carico degli utenti previsto. Questo processo prevede il provisioning di un ambiente cloud simile a quello di produzione, l'utilizzo di strumenti di test di carico per generare il carico e l'analisi dei parametri per valutare la capacità del carico di lavoro di gestire un carico realistico. Occorre eseguire i test di carico tramite versioni sintetiche o purificate dei dati di produzione (rimuovendo le informazioni sensibili o che permettono l'identificazione degli utenti). Eseguite automaticamente i test di carico come parte della vostra pipeline di distribuzione e confrontate i risultati con soglie e soglie predefinite KPIs. Questo processo ti consente di ottenere le prestazioni richieste.

Passaggi dell'implementazione

- Definisci gli obiettivi dei test: individua gli aspetti in termini di prestazione del carico di lavoro da valutare, come il throughput e il tempo di risposta.
- Seleziona uno strumento di test: scegli e configura lo strumento di test più adatto al carico di lavoro.
- Configura l'ambiente: configura l'ambiente di test in base al tuo ambiente di produzione. Puoi utilizzare AWS i servizi per eseguire ambienti su scala di produzione per testare la tua architettura.

- Implementa il monitoraggio: utilizza strumenti di monitoraggio come [Amazon CloudWatch](#) per raccogliere metriche tra le risorse della tua architettura. Puoi anche raccogliere e pubblicare metriche personalizzate.
- Definisci gli scenari definisci scenari e parametri del test di carico (come la durata del test e il numero di utenti).
- Esegui test di carico: effettua scenari di test su vasta scala. Approfittane Cloud AWS per testare il tuo carico di lavoro e scoprire dove non riesce a scalare o se è scalabile in modo non lineare. Ad esempio, usa le istanze spot per generare carichi a costi ridotti e rilevare i colli di bottiglia prima che si verifichino in produzione.
- Analizza i risultati dei test: analizza i risultati per individuare colli di bottiglia delle prestazioni e aree di miglioramento.
- Documenta e condividi gli esiti: documenta esiti e raccomandazioni e crea report al riguardo. Condividi queste informazioni con le parti interessate per aiutarle a prendere decisioni informate sulle strategie di ottimizzazione delle prestazioni.
- Effettua iterazioni continue: esegui con regolarità i test di carico, specie dopo una modifica o un aggiornamento del sistema.

Risorse

Documenti correlati:

- [Amazon CloudWatch RUM](#)
- [Amazon CloudWatch Synthetics](#)
- [Test di carico distribuito su AWS](#)

Video correlati:

- [AWS Summit ANZ 2023: accelera con fiducia grazie ai test di carico AWS distribuiti](#)
- [AWS re:Invent 2022: scalabile AWS per i primi 10 milioni di utenti](#)
- [Soluzione con AWS soluzioni: test di carico distribuiti](#)
- [AWS re:Invent 2021 - Ottimizza le applicazioni attraverso approfondimenti sugli utenti finali con Amazon CloudWatch RUM](#)
- [Demo di Amazon CloudWatch Synthetics](#)

Esempi correlati:

- [Test di carico distribuito su AWS](#)

PERF05-BP05 Usa l'automazione per risolvere in modo proattivo i problemi relativi alle prestazioni

Utilizzate gli indicatori chiave di prestazione (KPIs), combinati con i sistemi di monitoraggio e avviso, per affrontare in modo proattivo i problemi relativi alle prestazioni.

Anti-pattern comuni:

- Solo il personale operativo è autorizzato ad apportare modifiche operative al carico di lavoro.
- Tutti gli allarmi giungono direttamente al team operativo senza alcuna correzione proattiva.

Vantaggi dell'adozione di questa best practice: la correzione proattiva delle azioni di allarme consente al personale di supporto di concentrarsi sugli elementi non attivabili in automatico. In questo modo, il personale operativo non viene sovraccaricato da tutti gli allarmi e si concentra, invece, solo sugli allarmi critici.

Livello di rischio associato se questa best practice non fosse adottata: basso

Guida all'implementazione

Laddove possibile, utilizza gli allarmi per attivare operazioni automatizzate per risolvere i problemi. Se non è possibile rispondere in modo automatizzato, inoltra l'allarme a chi può intervenire. Ad esempio, potreste disporre di un sistema in grado di prevedere i valori previsti degli indicatori chiave di prestazione (KPI) e di avvisare quando superano determinate soglie, oppure uno strumento in grado di arrestare o ripristinare automaticamente le implementazioni se non raggiungono i valori previsti. KPIs

Implementa processi che forniscono visibilità sulle prestazioni durante l'esecuzione del carico di lavoro. Crea pannelli di controllo del monitoraggio e stabilisci norme di riferimento per le aspettative in termini di prestazioni, per determinare se il carico di lavoro presenta prestazioni ottimali.

Passaggi dell'implementazione

- Identifica il flusso di correzione: individua e comprendi il problema delle prestazioni risolvibile automaticamente. Utilizza soluzioni di AWS monitoraggio come [Amazon CloudWatch](#) o AWS X-Ray per aiutarti a comprendere meglio la causa principale del problema.
- Definisci il processo di automazione: crea un processo di step-by-step riparazione che può essere utilizzato per risolvere automaticamente il problema.
- Configura l'evento di avvio: configura l'evento per l'avvio automatico del processo di risoluzione. Ad esempio, è possibile definire un trigger per riavviare automaticamente un'istanza quando raggiunge una determinata soglia di CPU utilizzo.
- Automatizza la riparazione: utilizza AWS servizi e tecnologie per automatizzare il processo di riparazione. Ad esempio, [AWS Systems Manager Automation](#) fornisce un modo sicuro e scalabile per automatizzare il processo di risoluzione. Assicurati di utilizzare la logica di risoluzione automatica per annullare le modifiche se non risolvono correttamente il problema.
- Testa il flusso di lavoro: esegui il test del processo di risoluzione automatizzato in un ambiente di preproduzione.
- Implementa il flusso di lavoro: implementa la risoluzione automatizzata nell'ambiente di produzione.
- Sviluppa un playbook: predisponi e documenta un playbook che delinei le fasi del piano di risoluzione, inclusi eventi di avvio, logica di risoluzione e azioni intraprese. Assicurati di fornire la giusta preparazione alle parti interessate in modo che possano rispondere efficacemente agli eventi di risoluzione automatizzati.
- Esamina e perfeziona: valuta con regolarità l'efficacia del flusso di lavoro di risoluzione automatizzato. Modifica gli eventi di avvio e la logica di risoluzione, se necessario.

Risorse

Documenti correlati:

- [CloudWatchDocumentazione](#)
- [Partner per il monitoraggio, la registrazione e le prestazioni AWS Partner Network](#)
- [Documentazione di X-Ray](#)
- [Utilizzo di allarmi e azioni di allarme in CloudWatch](#)
- [Sviluppa una pratica di automazione del cloud per l'eccellenza operativa: le migliori pratiche di AWS Managed Services](#)

- [Automate your Amazon Redshift performance tuning with automatic table optimization](#)

Video correlati:

- [AWS re:Invent 2023 - Strategie per la scalabilità automatizzata, la correzione e l'autoguarigione intelligente](#)
- [AWS re:Invent 2023 - \[\] Monitoraggio delle applicazioni per carichi di lavoro moderni LAUNCH](#)
- [AWS re:Invent 2023 - Implementazione dell'osservabilità delle applicazioni](#)
- [AWS re:Invent 2021 - Automatizzazione intelligente delle operazioni cloud](#)
- [AWS re:Invent 2022 - Configurazione di controlli su larga scala nel proprio ambiente AWS](#)
- [AWS re:Invent 2022 - Automatizzazione della gestione e della conformità delle patch utilizzando AWS](#)
- [AWS re:Invent 2022 - In che modo Amazon utilizza metriche migliori per migliorare le prestazioni del sito Web](#)
- [AWS re:Invent 2023 - Prenditi una pausa: diagnostica e risolvi i problemi di prestazioni con Amazon RDS](#)
- [AWS re:Invent 2021 - {New Launch} Rileva e risolvi automaticamente i problemi con Amazon Guru DevOps](#)
- [AWS re:Invent 2023 - Centralizza le tue operazioni](#)

Esempi correlati:

- [CloudWatch Registri: personalizza gli allarmi](#)

PERF05-BP06 Conserva il carico di lavoro e i servizi up-to-date

Resta up-to-date su nuovi servizi e funzionalità cloud per adottare funzionalità efficienti, rimuovere problemi e migliorare l'efficienza complessiva delle prestazioni del tuo carico di lavoro.

Anti-pattern comuni:

- Si ritiene che l'architettura corrente diventi statica e non venga aggiornata nel corso del tempo.
- Non si dispone di sistemi né si esegue regolarmente una valutazione per la compatibilità di software e pacchetti aggiornati con il carico di lavoro.

Vantaggi derivanti dall'adozione di questa best practice: stabilendo un processo per rimanere aggiornato up-to-date su nuovi servizi e offerte, puoi adottare nuove funzionalità e funzionalità, risolvere problemi e migliorare le prestazioni del carico di lavoro.

Livello di rischio associato se questa best practice non fosse adottata: basso

Guida all'implementazione

Valuta i modi per migliorare le prestazioni man mano che nuovi servizi, modelli di progettazione e funzionalità di prodotti diventano disponibili. Determina in che modo possono migliorare le prestazioni o aumentare l'efficienza del carico di lavoro tramite valutazioni, discussioni interne o analisi esterne. Definisci un processo per valutare gli aggiornamenti, le nuove funzionalità e i servizi pertinenti per il tuo carico di lavoro. Ad esempio, crea un proof of concept che utilizza le nuove tecnologie o consultati con un gruppo interno. Quando provi nuove idee o servizi, esegui test delle prestazioni per misurare l'impatto sulle prestazioni del carico di lavoro.

Passaggi dell'implementazione

- Esegui l'inventario del tuo carico di lavoro: esegui l'inventario di software e architettura del carico di lavoro e identifica i componenti da aggiornare.
- Identifica le origini dell'aggiornamento: identifica novità e origini dell'aggiornamento relative ai componenti del carico di lavoro. Ad esempio, puoi iscriverti al [AWS blog What's New at](#) per i prodotti che corrispondono al tuo componente di carico di lavoro. Puoi iscriverti al RSS feed o gestire le tue [iscrizioni e-mail](#).
- Definisci un programma di aggiornamento: definisci un programma per valutare nuovi servizi e funzionalità per il tuo carico di lavoro.
 - Puoi utilizzare [AWS Systems Manager Inventory](#) per raccogliere i metadati del sistema operativo (OS), delle applicazioni e delle istanze dalle tue EC2 istanze Amazon e capire rapidamente quali istanze eseguono il software e le configurazioni richieste dalla tua politica software e quali istanze devono essere aggiornate.
- Valuta il nuovo aggiornamento: individua le modalità di aggiornamento dei componenti del carico di lavoro. Sfrutta l'agilità del cloud per testare in modo semplice e rapido il modo in cui le nuove funzionalità possono migliorare il carico di lavoro per ottenere efficienza delle prestazioni.
- Utilizza l'automazione: sfrutta l'automazione del processo di aggiornamento per ridurre il livello di impegno per implementare le nuove funzionalità e limitare gli errori causati dai processi manuali.
 - Puoi utilizzare [CI/CD](#) per aggiornare AMLs automaticamente le immagini dei container e altri elementi relativi alla tua applicazione cloud.

- È possibile utilizzare strumenti come [AWS Systems Manager Patch Manager](#) per automatizzare il processo di aggiornamento del sistema e pianificare l'attività utilizzando le [finestre di manutenzione di AWS Systems Manager](#).
- Documenta il processo: documenta il tuo processo di valutazione di aggiornamenti e nuovi servizi. Fornisci ai proprietari il tempo e lo spazio necessari per ricercare, testare, sperimentare e convalidare aggiornamenti e nuovi servizi. Fate riferimento ai requisiti aziendali documentati e aiutateci KPIs a stabilire le priorità degli aggiornamenti che avranno un impatto positivo sull'azienda.

Risorse

Documenti correlati:

- [Blog AWS](#)
- [Cosa c'è di nuovo con AWS](#)
- [Implementazione di up-to-date immagini con pipeline automatizzate di EC2 Image Builder](#)

Video correlati:

- [AWS RE:InForce 2022 - Automattizzazione della gestione e della conformità delle patch utilizzando AWS](#)
- [All Things Patch: | Eventi AWS Systems ManagerAWS](#)

Esempi correlati:

- [Inventory and Patch Management](#)
- [One Observability Workshop](#)

PERF05-BP07 Rivedi le metriche a intervalli regolari

Nell'ambito della manutenzione ordinaria o in risposta a eventi o incidenti, esamina i parametri raccolti. Stabilisci quali di questi parametri sono fondamentali per risolvere i problemi e quali altri parametri aggiuntivi, se monitorati, possono contribuire a identificare, affrontare o prevenire i problemi.

Anti-pattern comuni:

- Si lascia che i parametri rimangano in uno stato di allarme per un lungo periodo di tempo.
- Creazione di allarmi non utilizzabili da un sistema di automazione.

Vantaggi dell'adozione di questa best practice: esamina in modo continuo i parametri raccolti per verificare che identifichino, risolvano o prevengano adeguatamente i problemi. I parametri possono anche diventare obsoleti se lasciati in uno stato di allarme per un lungo periodo di tempo.

Livello di rischio associato se questa best practice non fosse adottata: medio

Guida all'implementazione

Migliora continuamente la raccolta e il monitoraggio dei parametri. Nell'ambito della risposta a incidenti ed eventi, valuta quali parametri sono stati utili per affrontare il problema e quali sarebbero stati utili ma non sono attualmente misurati. Questo metodo ti aiuterà a migliorare la qualità dei parametri raccolti, in modo da prevenire o risolvere in modo più rapido gli incidenti futuri.

Nell'ambito della risposta a incidenti ed eventi, valuta quali parametri sono stati utili per affrontare il problema e quali sarebbero stati utili ma non sono attualmente misurati. Queste considerazioni ti aiuteranno a migliorare la qualità dei parametri raccolti, così da prevenire o risolvere più rapidamente gli incidenti futuri.

Passaggi dell'implementazione

- Definisci metriche: stabilisci metriche in termini di prestazioni critiche da monitorare, allineate all'obiettivo del carico di lavoro, incluse metriche quali il tempo di risposta e l'utilizzo delle risorse.
- Stabilisci una base: imposta un valore di base e auspicabile per ciascuna metrica. La base deve fornire i punti di riferimento per identificare deviazioni o anomalie.
- Imposta una cadenza: imposta una cadenza (ad esempio, settimanale o mensile) per rivedere le metriche più critiche.
- Identifica i problemi di prestazioni: durante ogni revisione, valuta tendenze e deviazione dai valori di base. Cerca eventuali colli di bottiglia o anomalie nelle prestazioni. Per i problemi identificati, esegui un'analisi approfondita delle cause principali per comprendere il motivo più importante alla base del problema.
- Individua le azioni correttive: utilizza l'analisi per identificare le azioni correttive, come l'ottimizzazione dei parametri, la correzione di bug e il dimensionamento delle risorse.
- Documenta gli esiti: documenta gli esiti, compresi i problemi identificati, le cause principali e le azioni correttive.

- Itera migliora: valuta e migliora continuamente il processo di revisione delle metriche. Usa le indicazioni apprese dalla revisione precedente per migliorare il processo nel tempo.

Risorse

Documenti correlati:

- [CloudWatch Documentazione](#)
- [Raccogli parametri e log da Amazon EC2 Instances e server locali con l'agente CloudWatch](#)
- [Interroga le tue metriche con Metrics Insights CloudWatch](#)
- [Partner per il monitoraggio, la registrazione e le prestazioni AWS Partner Network](#)
- [Documentazione di X-Ray](#)

Video correlati:

- [AWS re:Invent 2022 - Configurazione di controlli su larga scala nel proprio ambiente AWS](#)
- [AWS re:Invent 2022 - In che modo Amazon utilizza metriche migliori per migliorare le prestazioni del sito Web](#)
- [AWS re:Invent 2023 - Creazione di un'efficace strategia di osservabilità](#)
- [AWS Summit SF 2022 - Osservabilità completa e monitoraggio delle applicazioni con AWS](#)
- [AWS re:Invent 2023 - Prenditi una pausa: diagnostica e risolvi i problemi di prestazioni con Amazon RDS](#)

Esempi correlati:

- [Creazione di una dashboard con Amazon QuickSight](#)
- [CloudWatch Pannelli di controllo](#)

Conclusioni

Raggiungere e mantenere l'efficienza delle prestazioni richiede un approccio basato sui dati. Devi prendere in considerazione in modo attivo gli schemi di accesso e i compromessi che ti permetteranno di ottimizzare ulteriormente le prestazioni. I processi di revisione basati su benchmark e test di carico ti permettono di selezionare i tipi di risorse e le configurazioni più adatte. Trattare l'infrastruttura come codice ti aiuta a fare evolvere l'architettura in modo rapido e sicuro, mentre potrai utilizzare i dati per prendere decisioni informate in merito all'architettura stessa. Adoperare una combinazione di monitoraggio attivo e passivo ti aiuterà a mantenere costanti le prestazioni dell'architettura nel corso del tempo.

AWS si impegna ad aiutarvi a creare architetture che funzionino in modo efficiente offrendo al contempo valore aziendale. Utilizza gli strumenti e le tecniche illustrati in questo documento per avere successo.

Collaboratori

Le seguenti persone e organizzazioni hanno contribuito a questo documento:

- Sam Mokhtari, Senior Efficiency Lead Solutions Architect, Amazon Web Services
- Josh Hart, Solutions Architect, Amazon Web Services
- Richard Trabing, Solutions Architect, Amazon Web Services
- Brett Looney, Principal Solutions Architect, Amazon Web Services
- Nina Vogl, Principal Solutions Architect, Amazon Web Services
- Eric Pullen, Solutions Architect, Amazon Web Services
- Julien Lépine, Specialist SA Manager, Amazon Web Services
- Ronnen Slasky, Solutions Architect, Amazon Web Services

Approfondimenti

Per ulteriori informazioni, consulta le seguenti risorse:

- [Framework AWS Well-Architected](#)
- [AWS Architecture Center](#)

Revisioni del documento

Per ricevere una notifica sugli aggiornamenti del presente whitepaper, iscriviti al feed RSS.

Modifica	Descrizione	Data
Aggiornamento secondario alle best practice	PERF03-BP04 è stato aggiornato con nuovi suggerimenti sui servizi.	6 novembre 2024
Linee guida sulle best practice aggiornate	Diversi aggiornamenti di entità minore per tutto il pilastro.	27 giugno 2024
Aggiornamento importante e ristrutturazione	<p>Il pilastro è stato ristrutturato in modo da avere cinque aree di best practice (3 in meno rispetto a prima). Il contenuto è stato raggruppato nelle cinque aree e aggiornato.</p> <p>Le nuove aree di best practice sono selezione dell'architettura, calcolo e hardware, gestione dei dati, rete e distribuzione dei contenuti e processi e cultura.</p>	3 ottobre 2023
Aggiornamento secondario	Rimozione del linguaggio non inclusivo.	13 aprile 2023
Aggiornamenti per il nuovo framework	Best practice aggiornate con prontuario e nuove best practice aggiunte.	10 aprile 2023
Aggiornamento del whitepaper	Best practice aggiornate con nuova guida all'implementazione.	15 dicembre 2022

Aggiornamento del whitepaper	Ampliamento delle best practice e aggiunta dei piani di miglioramento.	20 ottobre 2022
Aggiornamento secondario	Rimozione del linguaggio non inclusivo.	22 aprile 2022
Aggiornamenti minori	Link aggiornati.	10 marzo 2021
Aggiornamenti minori	Timeout AWS Lambda modificato in 900 secondi e corretto il nome di Amazon Keyspaces (per Apache Cassandra).	5 ottobre 2020
Aggiornamento secondario	Correzione di un link danneggiato.	15 luglio 2020
Aggiornamenti per il nuovo framework	Revisione e aggiornamento importanti dei contenuti	8 luglio 2020
Aggiornamento del whitepaper	Aggiornamento minore per la correzione di problemi grammaticali	1° luglio 2018
Aggiornamento del whitepaper	Aggiornamento del whitepaper per rispecchiare le modifiche apportate a AWS	1° novembre 2017
Pubblicazione iniziale	Pubblicazione del pilastro dell'efficienza delle prestazioni - Framework AWS Well-Architected.	1° novembre 2016

Note

I clienti sono responsabili della propria valutazione indipendente delle informazioni contenute nel presente documento. Questo documento: (a) è solo a scopo informativo, (b) rappresenta le offerte e le pratiche attuali di AWS prodotti, che sono soggette a modifiche senza preavviso, e (c) non crea alcun impegno o assicurazione da parte dei suoi affiliati, AWS fornitori o licenzianti. AWS i prodotti o i servizi sono forniti «così come sono» senza garanzie, dichiarazioni o condizioni di alcun tipo, esplicite o implicite. Le responsabilità e le responsabilità dei AWS propri clienti sono regolate da AWS accordi e il presente documento non fa parte di, né modifica, alcun accordo tra AWS e i suoi clienti.

© 2023, Amazon Web Services, Inc. o società affiliate. Tutti i diritti riservati.

AWS Glossario

Per la AWS terminologia più recente, consultate il [AWS glossario](#) nella sezione Reference. Glossario AWS