



UNIVERSITÀ POLITECNICA DELLE MARCHE
FACOLTÀ DI INGEGNERIA

Corso di Laurea in Ingegneria Informatica e dell'Automazione

**STATO DELL'ARTE DELLE ARCHITETTURE HARDWARE DEDICATE
PER L'INTELLIGENZA ARTIFICIALE**

**STATE OF THE ART OF HARDWARE ARCHITECTURES DEDICATED
FOR ARTIFICIAL INTELLIGENCE**

Relatore: Chiar.ma

Prof. Diamantini Claudia

Tesi di Laurea di:

Coloccioni Jacopo

Correlatore: Chiar.mo

Prof. Bordi Nazzeno

A.A. 2022 / 2023

L'ottimizzazione delle architetture nell'implementazione dell'intelligenza artificiale (IA) è una sfida cruciale in un contesto in cui modelli di machine learning e deep learning, caratterizzati da milioni o miliardi di parametri, richiedono una gestione efficiente della complessità computazionale. La varietà delle risorse hardware, dalle CPU tradizionali alle GPU specializzate e alle Tensor Processing Units (TPU), ha amplificato la complessità di questa sfida.

Questo lavoro di tesi si propone di esplorare e analizzare le attuali architetture disponibili sia nei servizi cloud sia nei dispositivi embedded, focalizzandosi sull'implementazione dell'IA. L'obiettivo primario è identificare le architetture ottimizzate per supportare sia l'addestramento che l'inferenza su sistemi Cloud e dispositivi Edge, tenendo conto delle diverse configurazioni richieste da ciascuna piattaforma.

Nella sezione introduttiva, si sottolinea l'importanza critica di un'architettura efficiente e ottimizzata, necessaria per eseguire programmi di IA su sistemi cloud e dispositivi Edge. Si evidenzia la necessità di adattare le architetture alle specifiche esigenze di addestramento e inferenza in diversi contesti.

Il capitolo sulle architetture per l'IA inizia a delineare il corpo della tesi, esplorando i principali obiettivi dell'ottimizzazione, come la parallelizzazione massiccia e la gestione della larghezza di banda. Successivamente, vengono analizzate le architetture System-on-Chip (SoC) per l'IA, con un focus particolare sull'Apple A17 Pro Bionic.

Un intero capitolo è dedicato agli acceleratori hardware progettati per migliorare le performance, come il chip GAUDI2 di Intel per l'addestramento su cloud e l'AWS Inferentia2 per l'inferenza cloud.

La tesi approfondisce ulteriormente il ruolo delle GPU, nonostante il loro consumo energetico maggiore rispetto ad alcune soluzioni specializzate. Le GPU rimangono ampiamente utilizzate nei servizi cloud per attività di addestramento e inferenza grazie alle loro prestazioni e scalabilità su larga scala.

Dopo le GPU andremo a descrivere l'architettura delle TPU, utilizzate da Google per l'addestramento delle loro IA.

Infine, sarà riportata una tabella riassuntiva per evidenziare le differenze dei vari chip e un esempio di architettura per l'implementazione di un'auto a guida autonoma di livello cinque facendo uso delle nozioni riportate nei capitoli precedenti.

Attraverso questa analisi esaustiva, la tesi mira a fornire una panoramica completa delle architetture attuali, sottolineando le sfide e le soluzioni nel campo dell'ottimizzazione per l'implementazione dell'IA.

Sommario

INTRODUZIONE.....	7
CHE COS'E' L'INTELLIGENZA ARTIFICIALE	7
ALGORITMI DI APPRENDIMENTO DELLE AI.....	9
EDGE COMPUTING vs CLOUD COMPUTING.....	13
MOTIVAZIONI SULLO SVILUPPO DELL'HARDWARE AI	16
ARCHITETTURE IA	20
SOC IA.....	22
Architetture di un SoC ottimizzato per l'IA e implementato su Smartphone.....	24
ACCELERATORI HARDWARE.....	26
Architettura di un Acceleratore ottimizzato per l'addestramento su Cloud	28
Architettura di un Acceleratore ottimizzato per l'inferenza su Cloud	35
Architettura delle GPU	37
Architettura di un'istanza che utilizza GPU per l'addestramento su Cloud	48
NPU	50
Architettura delle TPU	51
CONFRONTO TECNICO DELLE ARCHITETTURE OTTIMIZZATE PER L'IA	55
ESEMPIO DI UN'ARCHITETTURA PER L'IMPLEMENTAZIONE DI UN'AUTO A GUIDA AUTONOMA.....	57
STORAGE PER ADDESTRAMENTO DI MACCHINA A GUIDA AUTONOMA	59
RISORSE COMPUTAZIONALI PER ADDESTRAMENTO DI UNA IA PER UNA MACCHINA A GUIDA AUTONOMA.....	61
ARCHITETTURA DELLA MACCHINA A GUIDA AUTONOMA.....	64
MEMORIA DEL COMPUTER DI BORDO DI UN' AUTO A GUIDA AUTONOMA.....	66
UNITA' DI ELABORAZIONE DEL COMPUTER DI BORDO DI UN' AUTO A GUIDA AUTONOMA.....	68
CONCLUSIONI	73

RINGRAZIAMENTI	74
FONTI BIBLIOGRAFICHE E SITOGRAFIA	75

INTRODUZIONE

CHE COS'E' L'INTELLIGENZA ARTIFICIALE

Spiegazione generale sull'Intelligenza Artificiale e Applicazioni Moderne

L'intelligenza artificiale (IA) è un ramo dell'informatica che si occupa della di creare macchine intelligenti in grado di eseguire compiti che tipicamente richiedono l'intelligenza umana attraverso lo sviluppo di algoritmi [77]; sebbene essa possa sembrare una novità degli ultimi anni, il concetto di ricreare un calcolatore che si comporti come un essere senziente ha affascinato l'idea di molti scienziati e informatici da un secolo a questa parte, basti pensare che già nel 1950 Alan Turing pubblicò un articolo chiamato "Computing Machinery and Intelligence" in cui si chiedeva se le macchine potevano pensare, offrendo tra l'altro un test definito "Test di Turing ¹" in cui un interrogatore umano proverebbe a distinguere tra una risposta testuale del computer e una umana.

Non fu però Turing a coniare il termine Intelligenza Artificiale ma un professore di matematica, John McCarthy, nel 1956, durante uno storico seminario che ruotava intorno alla possibilità di far realizzare alle macchine imprese dove era necessaria l'intelligenza umana.

Inizialmente, la ricerca sull'Intelligenza Artificiale si è concentrata principalmente sullo sviluppo di algoritmi. Tuttavia, l'introduzione dei microprocessori nel 1971, con l'Intel 4004, ha aperto la strada alla creazione di vere e proprie IA. Un esempio è rappresentato da R1, utilizzato dalla Digital Equipment nel 1982, il cui scopo era assistere nella configurazione degli ordini per nuovi computer. R1 fu la prima di molte IA utilizzate in ambito commerciale.

L'evoluzione successiva dei microprocessori, guidata dalla legge di Moore e dalla tecnologia MOS, ha notevolmente migliorato le prestazioni e l'efficienza dei calcolatori [1] [2]. Inoltre, grazie all'avanzamento delle tecnologie di memorizzazione, come le nuove RAM DDR e le unità SSD, e alla creazione di nuove architetture di bus come il PCIe o il massiccio utilizzo della parallelizzazione, è stato possibile sviluppare IA sempre più complesse e rivoluzionarie.

¹ Il Test di Turing è un criterio per valutare l'intelligenza artificiale basandosi sulla capacità di una macchina di emulare comportamenti indistinguibili da quelli umani. Tuttavia, è stato oggetto di critiche per le sue definizioni ambigue di "intelligenza" e la possibilità di superare il test attraverso emulazioni superficiali.

Attualmente, l'intelligenza artificiale ha raggiunto uno sviluppo significativo, integrandosi in una vasta gamma di settori e diventando essenziale per la nostra quotidianità; basti pensare agli assistenti vocali come Amazon Alexa o Siri, oppure ai software per migliorare le fotografie sui nostri smartphone.

Possiamo trovare questa tecnologia anche in molti altri esempi, come in Google Maps per determinare il miglior percorso o nelle piattaforme social per suggerire agli utenti gli account da seguire, ma si possono fare centinaia di esempi ben più complessi come le auto a guida autonoma di Tesla o le IA generative come ChatGPT di OpenAI.

ALGORITMI DI APPRENDIMENTO DELLE AI

Spiegazione delle Tipologie di Algoritmi di Apprendimento concentrando l'Attenzione sulla loro Complessità

Quando parliamo di IA dobbiamo considerare innanzitutto che ne esistono più tipologie, ognuna delle quali sarà contraddistinta da differenti algoritmi, ma in particolare ci focalizzeremo su una metodologia fondamentale che sottende all'efficacia dell'Intelligenza Artificiale: il Machine Learning (ML).

Machine Learning

Il Machine Learning rappresenta la capacità dei computer di apprendere dall'esperienza, anche se questo termine deve essere inteso in un contesto più ristretto rispetto all'apprendimento umano. In questo contesto informatico, l'apprendimento avviene quando le prestazioni di un programma migliorano dopo l'esecuzione di un compito o il completamento di un'azione, anche se l'azione può essere errata, seguendo il principio che anche gli errori possono contribuire all'apprendimento.

Nel contesto del Machine Learning, non si tratta di scrivere passo dopo passo il codice di programmazione che dice alla macchina cosa fare. Invece, al programma vengono forniti insiemi di dati, e attraverso l'applicazione di specifici algoritmi, la macchina sviluppa autonomamente una logica per eseguire una funzione orientata a uno specifico obiettivo. In altre parole, la macchina "impara" dai dati e adatta il proprio comportamento in base a questi input, migliorando le prestazioni nel tempo senza che sia necessario specificare esplicitamente ogni dettaglio del processo attraverso il codice. Questo approccio consente di affrontare complessità e variabilità nei dati senza richiedere una programmazione dettagliata per ogni possibile scenario [17] [23].

Esistono poi numerosi sottoinsiemi che consentono di classificare in modo più dettagliato il Machine Learning in base al suo funzionamento:

Apprendimento Supervisionato

L'apprendimento supervisionato implica l'addestramento degli algoritmi mediante l'utilizzo di esempi già classificati, in cui gli input sono associati agli output corrispondenti. Ad esempio, consideriamo

un dataset in cui i dati sono classificati come "F" (failed) o "R" (runs). Durante il processo di addestramento, l'algoritmo apprende associando input con gli output.

Questo apprendimento avviene attraverso la comparazione dei risultati ottenuti con quelli attesi. Gli errori sono identificati e il modello viene modificato di conseguenza per migliorare la sua capacità predittiva. L'apprendimento supervisionato si avvale di tecniche come la classificazione, la regressione, la previsione e il gradient boosting. Tali modelli sono in grado di predire valori per dati non ancora classificati.

Un esempio concreto di utilizzo dell'apprendimento supervisionato è la prevenzione delle frodi nelle transazioni con carte di credito o la previsione delle richieste di risarcimento da parte di clienti in un'azienda assicurativa. In questi contesti, l'apprendimento supervisionato sfrutta i dati storici per anticipare eventi futuri e prendere decisioni informate.

Apprendimento Semi Supervisionato

L'apprendimento semi-supervisionato condivide molte delle stesse applicazioni dell'apprendimento supervisionato, ma presenta alcune caratteristiche distintive. In questo approccio, l'addestramento dell'algoritmo coinvolge sia dati classificati che dati non classificati. Di solito, si dispone di un volume limitato di dati classificati e di un volume più ampio di dati non classificati, poiché acquisire quest'ultimo è spesso più economico e meno laborioso.

Durante il processo di addestramento, l'algoritmo utilizza i dati classificati per apprendere i pattern e le relazioni tra gli input e gli output noti, ma allo stesso tempo, sfrutta i dati non classificati per espandere la sua comprensione e adattarsi a una varietà più ampia di situazioni.

L'apprendimento semi-supervisionato può essere applicato attraverso metodi di classificazione, regressione e previsione. Questo approccio è particolarmente vantaggioso quando il processo di classificazione ha un costo elevato e non è pratico o sostenibile eticamente classificare manualmente tutti i dati.

Un esempio recente di successo nell'apprendimento semi-supervisionato è rappresentato dalle fotocamere in grado di identificare i volti delle persone. In questo caso, l'algoritmo può essere addestrato su un insieme di dati contenente volti classificati, ma può anche generalizzare ulteriormente la sua capacità di riconoscere volti sfruttando un vasto numero di dati non classificati. Questo approccio rende l'apprendimento più scalabile ed efficiente.

Apprendimento Non Supervisionato

L'apprendimento non supervisionato è un approccio utilizzato su dati che non sono pre-classificati. In questo contesto, il sistema non riceve una "risposta giusta"; piuttosto, l'algoritmo è incaricato di scoprire autonomamente la struttura interna dei dati presentati. L'obiettivo principale è esplorare i dati senza fornire indicazioni esplicite, individuando eventuali pattern o relazioni intrinseche.

Questo tipo di apprendimento risulta particolarmente efficace con dati transazionali. Ad esempio, può identificare consumatori con caratteristiche simili, permettendo di indirizzare campagne di marketing specifiche a gruppi omogenei. In alternativa, può individuare le caratteristiche chiave che differenziano segmenti di consumatori l'uno dall'altro.

Questi algoritmi sono impiegati in diverse applicazioni, come la segmentazione di argomenti testuali, la raccomandazione di prodotti o l'identificazione di valori anomali nei dati. In sintesi, l'apprendimento non supervisionato consente di esplorare e comprendere la struttura interna di dati complessi senza la necessità di etichettare in anticipo le informazioni contenute.

Apprendimento per Rinforzo

L'apprendimento per rinforzo è spesso impiegato in settori come la robotica, i videogiochi e la navigazione, offrendo un approccio unico all'apprendimento in base all'esperienza e alle ricompense. In questo contesto, l'algoritmo impara attraverso esperimenti ed errori, scoprendo quali azioni generano le ricompense più elevate.

Questo tipo di apprendimento presenta tre componenti principali: l'agente (responsabile dell'apprendimento o delle decisioni), l'ambiente (l'insieme di tutto con cui l'agente interagisce) e le azioni (le possibili azioni che l'agente può intraprendere). L'obiettivo fondamentale dell'agente è selezionare le azioni che massimizzano le ricompense previste in un determinato periodo temporale. Facendo le scelte corrette, l'agente può raggiungere l'obiettivo in modo più efficiente.

Quindi, l'obiettivo finale dell'apprendimento per rinforzo è acquisire la conoscenza necessaria per determinare le azioni ottimali da intraprendere. Questo approccio è particolarmente utile in situazioni in cui l'agente deve apprendere a navigare in un ambiente dinamico, risolvere problemi complessi o ottimizzare le sue azioni per massimizzare le ricompense nel lungo termine.

Deep Learning

Il Deep Learning è una sotto-disciplina del Machine Learning che si focalizza sull'uso di reti neurali artificiali profonde per apprendere automaticamente rappresentazioni complesse dei dati. Le reti neurali profonde sono composte da più strati (o livelli) di neuroni artificiali, e questa struttura stratificata consente di riconoscere caratteristiche sempre più complesse nei dati, rendendo il Deep Learning particolarmente adatto per compiti complessi come il riconoscimento di immagini e il trattamento del linguaggio naturale.

In particolare, questi neuroni sono moduli software chiamati nodi, che utilizzano calcoli matematici su degli input per elaborare i dati. Le reti neurali artificiali sono algoritmi di deep learning che utilizzano questi nodi per risolvere problemi complessi. [49]

ML o DL

La distinzione tra Machine Learning e Deep Learning risiede nella gestione dei dati e nel processo di apprendimento. Il ML è flessibile, trattando sia dati strutturati che non strutturati, con un coinvolgimento umano significativo nel processo di addestramento. Al contrario, il DL è spesso associato a dati non strutturati come immagini e testo, richiedendo un considerevole volume di dati e potenza di calcolo.

Mentre entrambi i campi incorporano l'autoapprendimento, il ML può richiedere più intervento umano nella configurazione e nella preparazione dei dati. D'altro canto, il DL coinvolge l'uomo nella progettazione dell'architettura della rete, anche se spesso richiede meno intervento durante il processo di addestramento.

In sostanza il DL è stato creato per sopperire alle mancanze e alle criticità degli algoritmi sopra citati del ML; tuttavia ci sono alcune difficoltà nella sua implementazione pratica, come ad esempio la necessità di una grande quantità di dati (maggiore rispetto a quelli del ML), oppure un insieme di dati di alta qualità, poiché l'utilizzo di valori anomali o presentanti degli errori potrebbero influire in modo negativo sull'apprendimento ed infine, in genere, per il DL sarà necessaria una elaborazione più intensiva e un'infrastruttura con capacità di calcolo sufficiente per funzionare correttamente altrimenti impiegheranno molto tempo per elaborare i risultati.

EDGE COMPUTING vs CLOUD COMPUTING

Differenze tra Edge e Cloud Computing per l'IA

Nei pochi esempi precedenti è possibile osservare anche una distinzione tra le intelligenze artificiali che si appoggiano a sistemi Cloud come ChatGPT e intelligenze artificiali che invece si basano completamente sul dispositivo di utilizzo come Siri.

Credo che sia quindi il caso di fare una differenziazione tra Edge Computing e Cloud Computing per quanto riguarda l'IA; infatti, questo sarà un tema centrale per lo sviluppo delle architetture in futuro.

Cloud Computing

Il Cloud Computing, con la sua architettura centralizzata basata su data center remoti, ha rivoluzionato la gestione dei dati aziendali. Il paradigma del cloud computing si fonda sull'idea che le operazioni di elaborazione principali avvengano su una macchina remota, distante da quella attualmente in uso dall'utente, che riceve direttamente e gestisce gli input [56]. I dati generati durante questo processo vengono archiviati ed elaborati da server remoti, noti anche come server cloud. Questo approccio significa che il dispositivo utilizzato per accedere al cloud non deve eseguire operazioni computazionali altrettanto intense.

Ospitando software, piattaforme e database in ambienti remoti, i server cloud liberano la memoria e la potenza di calcolo dei singoli dispositivi. Gli utenti possono in modo sicuro accedere ai servizi cloud mediante le credenziali fornite dal fornitore di servizi di cloud computing.

Poiché il cloud computing implica che il carico di lavoro del computer dell'utente è gestito su una macchina separata, l'accesso al cloud è possibile ovunque e aperto a chiunque disponga di una connessione Internet.

I data center inoltre hanno a disposizione strutture hardware molto potenti, con sofisticati sistemi di parallelizzazione, sistemi di raffreddamento avanzati e moltissimi dati a disposizione che gli permettono di essere l'ideale per moltissime applicazioni come lo sviluppo delle IA.

Tuttavia, ci sono sfide da affrontare. La connettività Internet affidabile è essenziale per i sistemi di machine learning basati su cloud, poiché la trasmissione dei dati grezzi al servizio cloud e il recupero dei dati elaborati dipendono da una connessione robusta. Anche se l'elaborazione nel cloud è più rapida rispetto a quella convenzionale, c'è un intervallo tra la trasmissione e la ricezione dei dati che

non è sempre accettabile, soprattutto quando si utilizzano algoritmi di apprendimento automatico, dove la velocità è fondamentale, oppure quando si utilizzano IA che richiedono grandi quantità di dati in input. È impossibile pensare ad esempio a macchine con guida autonoma che lavorino su Cloud, basterebbe infatti un guasto di rete per creare disagi considerevoli.

Ecco perché alcune applicazioni richiedono il supporto all'Edge, nel quale sicuramente vi sono architetture meno sofisticate e prestazionali, ma possono evitare alcuni dei problemi prima descritti.

Edge Computing

L'Edge computing, quindi, è una forma di elaborazione che viene eseguita in sede o in prossimità di una particolare origine dati per elaborare i dati critici localmente, per inviarli in seguito a una repository centrale, permettendo risposte in tempo reale, e il risparmio di banda, inviando al data center informazioni già elaborate e quindi di minori dimensioni [7]. Ciò è particolarmente vantaggioso per le applicazioni di IA dove il tempo di risposta è un aspetto critico o i dati che devono essere inviati al data center per essere elaborati sono di notevoli dimensioni [56].

Ne sono un esempio oltre alle macchine a guida autonoma, le IA che supportano gli operatori sanitari o gli assistenti vocali; queste applicazioni, infatti, necessitano di una risposta immediata evitando che i dati vengano spediti al cloud tramite la rete facendo sì che vengano elaborati direttamente in loco.

Tuttavia, poiché l'Edge computing, come già detto, si basa su dispositivi periferici di dimensioni ridotte e costruiti per applicazioni specifiche, mette a disposizione delle architetture decisamente meno prestazionali di quelle che troviamo nei data center, e anche per quanto riguarda l'addestramento delle AI nei dispositivi Edge si trovano decisamente meno dati.

Integrazione Edge-Cloud

L'integrazione tra Edge Computing e Cloud Computing emerge come una prospettiva futura promettente. Questa combinazione consente di massimizzare i benefici di entrambi gli approcci, superando le rispettive limitazioni. L'Edge AI può gestire decisioni locali, mentre i servizi cloud apprendono continuamente per migliorare i modelli AI. Questa sinergia offre rapidità ed efficienza nelle applicazioni e nei dispositivi IoT, mantenendo la capacità di archiviazione e potenza di elaborazione del cloud.

In particolare, al giorno d'oggi il cloud è il luogo ideale per la formazione (addestramento dell'AI) perché fornisce l'accesso a vasti archivi di dati da più server e a più informazioni di un'applicazione IA, inoltre, il cloud può ridurre le spese perché consente alle unità di elaborazione grafica (GPU) o altre architetture dedicate di addestrare più modelli di intelligenza artificiale. Poiché l'addestramento avviene in modo intermittente su ciascun modello, la capacità non rappresenta un problema.

Con l'inferenza invece, fase in cui gli algoritmi di intelligenza artificiale applicano le conoscenze precedentemente acquisite durante l'addestramento per generare risposte in tempo reale, si assiste a un notevole calo nella gestione dei dati, se confrontato con la fase di addestramento, questo permette di evitare una trasmissione di informazioni con i data center cloud e quindi di poter sfruttare tutti i vantaggi offerti dall'elaborazione sui dispositivi periferici [57].

MOTIVAZIONI SULLO SVILUPPO DELL'HARDWARE AI

Spiegazione sull'importanza e sullo sviluppo dell'hardware ottimizzato per l'IA

Come accennato in precedenza, le attuali intelligenze artificiali richiedono capacità computazionali complesse dei sistemi hardware e considerevoli quantità di dati su cui addestrarsi; pertanto, ottimizzare l'hardware e creare delle tecnologie specifiche permette di evitare colli di bottiglia rappresentati dalle limitazioni fisiche dei computer general-purpose² e in maniera tale da poter sviluppare IA sempre più sofisticate che possano essere un aiuto per l'uomo [9].

Complessità computazionale

Quando si parla di “complessità computazionale” ci si riferisce alla grande mole di operazioni matematiche simultanee, come moltiplicazioni di matrici e calcoli di gradienti, e quindi dell'algebra lineare, per ricavare un valore di output tramite un valore di input e una certa funzione che determina appunto, il risultato migliore.

L'importanza dell'Algebra Lineare nell'ambito dell'Intelligenza Artificiale emerge chiaramente negli algoritmi di ottimizzazione, una componente cruciale nell'addestramento dei modelli di apprendimento automatico.

Questi algoritmi si concentrano sull'obiettivo di minimizzare o massimizzare una funzione specifica, come l'errore di predizione di un modello, che viene ottenuta attraverso l'iterativo aggiornamento dei parametri del modello.

Il processo di ottimizzazione può essere concepito come un percorso in uno spazio vettoriale multidimensionale, dove ogni punto rappresenta un possibile set di parametri per il modello. In questo contesto, l'algebra lineare fornisce gli strumenti fondamentali per calcolare la direzione e la lunghezza del passo da compiere ad ogni iterazione.

Ad esempio, l'utilizzo di vettori gradienti consente di determinare la direzione in cui la funzione obiettivo cresce più rapidamente, indicando la strada per ridurre l'errore di predizione. Inoltre, il

² Un computer "general purpose" è un sistema di elaborazione dati progettato per eseguire una vasta gamma di applicazioni e compiti, a differenza di un sistema specializzato progettato per una funzione specifica. I computer general purpose sono in grado di eseguire diverse attività grazie al loro hardware e software flessibili, rendendoli adatti a scopi diversi, come elaborazione dati, calcoli complessi, gestione di informazioni, e altro ancora.

calcolo degli autovalori e degli autovettori può essere impiegato per determinare la lunghezza ottimale del passo da compiere in modo da raggiungere il minimo o il massimo della funzione obiettivo nel modo più efficiente possibile.

L'algebra lineare svolge un ruolo fondamentale, costituendo il fondamento matematico su cui si basano numerosi algoritmi di Machine Learning; ad esempio, infatti, i dati di input per gli algoritmi sono comunemente rappresentati come vettori o matrici, e le operazioni eseguite su questi dati, come la somma, la moltiplicazione o la trasformazione, sono tutte operazioni di Algebra.

Un caso di studio emblematico dell'applicazione dell'Algebra Lineare nel Machine Learning è l'algoritmo di regressione lineare. Questo approccio statistico è utilizzato per prevedere una variabile di output basandosi su una o più variabili di input, modellando la relazione tra di esse come una funzione lineare.

Nella regressione lineare, i parametri del modello sono spesso rappresentati come un vettore, e l'obiettivo è individuare il vettore di parametri che minimizza la differenza tra le previsioni del modello e i valori effettivi di output. Questo problema di ottimizzazione può essere risolto mediante tecniche di algebra lineare, come la minimizzazione dei quadrati, che consiste nel trovare il vettore di parametri che minimizza la somma dei quadrati delle differenze tra le previsioni del modello e i valori effettivi di output.

La regressione lineare è solo uno degli innumerevoli algoritmi di Machine Learning che fanno largo uso dell'Algebra Lineare [50]. Altri esempi comprendono reti neurali, macchine a vettori di supporto (SVM) e analisi dei componenti principali (PCA), solo per citarne alcuni.

Anche per quanto concerne il Deep Learning, l'algebra lineare rappresenta una componente indispensabile. Una rete neurale, infatti, è un modello ispirato al funzionamento del cervello umano, composto da diversi strati di nodi o neuroni che possono trasmettere informazioni tra di loro. Ogni nodo in uno strato riceve input dai nodi nello strato precedente, applica una funzione di attivazione e invia l'output ai nodi nello strato successivo.

Le informazioni trasmesse tra i nodi di una rete neurale sono rappresentate come vettori. Ad esempio, l'input di una rete neurale potrebbe essere un'immagine, rappresentata come un vettore di pixel. Anche le informazioni trasmesse tra i nodi in uno strato e quelli nello strato successivo sono rappresentate come vettori.

Le connessioni tra i nodi di una rete neurale sono descritte da una matrice chiamata "matrice dei pesi". Ogni elemento di questa matrice rappresenta il peso di una connessione tra due nodi. La matrice dei pesi viene moltiplicata per il vettore di input per produrre il vettore di output, in un processo noto

come propagazione in avanti. L'addestramento di una rete neurale implica la modifica dei pesi delle connessioni affinché l'output della rete si avvicini il più possibile all'output desiderato [50].

Questi modelli di DL possono estendersi su centinaia di livelli, ognuno composto da migliaia di nodi interconnessi, generando un numero impressionante di parametri. Ad esempio, una rete neurale avanzata, come quella utilizzata in modelli di lingua naturale o visione computerizzata, può facilmente raggiungere miliardi di parametri.

La natura intensiva delle operazioni all'interno di queste reti è evidente nelle molteplici moltiplicazioni di matrici coinvolte durante le fasi di addestramento e inferenza. In una singola iterazione attraverso la rete, decine di miliardi di operazioni di moltiplicazione di matrici possono essere eseguite, illustrando la vastità delle computazioni coinvolte. Proprio per questo, per l'utilizzo dell'IA, si ricercano architetture con un grande tasso di parallelizzazione e una grande larghezza di banda [9].

Tutto questo è servito allo sviluppo di hardware dedicato per le applicazioni di IA, sia nei dispositivi Edge che in quelli presenti nei data center, e ovviamente sarà anche la motivazione principale alla base degli sviluppi futuri.

Array Sistolico

Nel contesto delle reti neurali, l'esecuzione di numerosi calcoli è fondamentale, e gran parte di questi calcoli può essere semplificata nella forma di una singola operazione chiamata moltiplicazione-accumulo (MAC). Quando si tratta di eseguire queste operazioni su larga scala, ci si trova ad affrontare sfide di efficienza computazionale.

Le unità di elaborazione tradizionali, come le CPU, operano in modo scalare, gestendo le istruzioni uno alla volta. Sebbene quest'approccio sia adatto a scopi generici, per migliorare le prestazioni in contesti specifici, come il machine learning, è possibile specializzarsi ulteriormente.

Le GPU, ad esempio, operano su vettori 1D, consentendo di eseguire calcoli su un intero elenco di dati contemporaneamente. Tuttavia, questa modalità è più adatta a compiti che richiedono il ripetersi del medesimo calcolo su un insieme di dati.

Per ottimizzare ulteriormente le prestazioni nel contesto delle reti neurali, si può adottare un'architettura di tipo matriciale. Dato che i dati di una rete neurale sono organizzati in una matrice, costruire una "macchina a matrice" può portare a notevoli miglioramenti. Concentrandosi

principalmente sull'operazione di moltiplicazione-accumulo (MAC), è possibile dedicare la maggior parte delle risorse del chip a questa attività, ignorando in gran parte altre istruzioni.

L'approccio per ottenere tali prestazioni matriciali è attraverso un'architettura nota come "array sistolico". Questa configurazione impiega un algoritmo hardware che organizza le celle su un chip per eseguire la moltiplicazione di matrici in modo efficiente. Gli array sistolici consentono il calcolo di più operazioni in parallelo durante ogni accesso alla memoria, riducendo i picchi di consumo energetico associati agli accessi alla memoria. Il flusso di dati attraverso le celle segue un modello orchestrato, simile a una catena di montaggio, consentendo un'elaborazione continua dei dati in modo pipeline attraverso molte unità di elaborazione. In sostanza, un array sistolico rappresenta un insieme di celle interconnesse, ognuna specializzata nell'esecuzione di un'operazione specifica, e il flusso di dati segue un percorso ottimizzato per massimizzare l'efficienza computazionale.

Altre motivazioni

Un altro motivo per sviluppare hardware per l'IA ottimizzato, oltre all'accelerazione computazionale può essere la miglior gestione delle risorse energetiche; un tema molto odierno.

Sviluppare hardware che possano avere prestazioni migliori con un dispendio energetico inferiore è infatti molto importante sia per una questione economica delle aziende, per la salvaguardia del nostro pianeta, ma in particolar modo perché mantenere una temperatura bassa durante l'utilizzo di questi chip, permette di raggiungere delle prestazioni ottimali e un dispendio energetico ridotto nei dispositivi embedded si traduce anche in una migliore salvaguardia della batteria.

Ovviamente poi, un ultimo motivo è certamente quello economico; è stato infatti previsto da un rapporto, stilato nel dicembre del 2018 dalla McKinsey & Company, che entro il 2025, i semiconduttori legati all'intelligenza artificiale rappresenteranno il 20% di tutta la domanda mondiale, il che si traduce in circa 67 miliardi di dollari di entrate [1].

ARCHITETTURE IA

L'hardware per l'intelligenza artificiale rappresenta un elemento critico nell'ecosistema tecnologico moderno. La sua importanza deriva dal fatto che consente alle applicazioni di IA di eseguire complessi calcoli e operazioni di apprendimento in modo efficiente [9]. Questi sistemi hardware sono progettati per ottimizzare le prestazioni dei modelli di machine learning e deep learning, riducendo i tempi di addestramento, ma anche per migliorare l'inferenza e l'efficienza complessiva delle applicazioni IA [13]. L'architettura di queste componenti hardware, a differenza di quelle che non richiedono l'elaborazione di algoritmi di IA, è sviluppata e incentrata sul gestire compiti intensivi in termini di calcolo parallelo; il parallelismo, infatti, è uno degli aspetti cruciali richiesti dagli algoritmi di IA, poiché più un chip riesce ad elaborare calcoli matriciali velocemente più esso sarà performante e apprenderà in maniera rapida [14] [11]. Va infatti precisato che le moltiplicazioni matriciali sono delle operazioni molto parallelizzabili, è infatti possibile andare a parallelizzare una moltiplicazione matriciale senza particolari sforzi, questo fa sì che più core siano presenti all'interno di un chip, più queste moltiplicazioni verranno velocizzate [8].

Quando si sviluppano questi hardware poi, un'altra caratteristica molto importante è la larghezza di banda [15], fondamentale per consentire un flusso efficiente dei dati tra la memoria e l'unità di elaborazione; per larghezza di banda, si intende la quantità massima di dati che possono essere trasferiti da un punto all'altro in un dato intervallo di tempo ed è strettamente legata alla dimensione del bus e alla velocità di trasmissione, ossia la velocità con cui i dati possono essere inviati o ricevuti attraverso il bus.

È ovvio che avere a disposizione una grande larghezza di banda permette a questi chip di trasferire in meno passaggi i dati dalla memoria principale ai processori e viceversa, andando a ridurre esponenzialmente i tempi di elaborazione, assicurare dei trasferiti rapidi e migliorando le prestazioni complessive. In genere per aumentarla, si opta per l'espansione dei canali bus, ossia l'aumento del numero di fili conduttori dedicati al trasporto delle informazioni o aumentando la frequenza di lavoro, cioè il numero di cicli di clock in un secondo. Questo miglioramento può anche essere ottenuto implementando più canali di accesso alle memorie, come nel caso delle tecnologie DDR5, che tuttavia, può comportare una maggiore complessità del sistema e richiedere una gestione più avanzata della memoria per sfruttare appieno le potenzialità della tecnologia.

Nel perseguire l'obiettivo di massimizzare la larghezza di banda, è essenziale considerare attentamente l'ottimizzazione dell'intero sistema. Un aumento indiscriminato della larghezza di banda

potrebbe comportare sfide significative, come un consumo energetico eccessivo, una maggiore complessità nella gestione delle risorse e un potenziale impatto negativo sulla portabilità e sulle dimensioni del dispositivo.

Inoltre, nel processo di progettazione, è fondamentale considerare la natura dei carichi di lavoro previsti per il sistema. Ad esempio, se il dispositivo è destinato a compiti leggeri, un'elevata larghezza di banda potrebbe non tradursi in miglioramenti sostanziali delle prestazioni, ma comporterebbe comunque gli svantaggi precedentemente menzionati. La scelta del giusto equilibrio tra larghezza di banda, consumo energetico ed efficienza complessiva del sistema diventa quindi cruciale [14].

Un'altra caratteristica particolarmente ambita è la dimensione della memoria [8] [15], soprattutto nei dispositivi dedicati all'addestramento, poiché un set di dati più ampio per il processo di allenamento contribuisce all'accuratezza dell'intelligenza artificiale. Tuttavia, è cruciale considerare che, nonostante le dimensioni della memoria non influiscano negativamente sulle performance, è necessario valutare attentamente il rapporto costi-benefici. L'inclusione di una memoria più ampia in un chip comporta inevitabilmente un aumento dei costi, e se il chip non sfruttasse appieno questa memoria aggiuntiva, potrebbe risultare in un investimento non ottimale. Pertanto, è essenziale tenere in considerazione l'utilizzo effettivo del chip al fine di massimizzare l'efficacia della memoria implementata e garantire un bilanciamento appropriato tra prestazioni e costi.

SOC IA

Descrizione dei System on Chip per l'IA

Una particolare tipologia di hardware ottimizzato che in questi anni sta riscuotendo molto successo, sono i System on Chip (SoC) ottimizzati per l'IA, questi dispositivi integrano diverse componenti essenziali, come processori, unità di elaborazione grafica (GPU), memorie e acceleratori hardware specifici per l'IA, in un unico chip [29]. L'utilizzo di questa tecnologia può essere molto vantaggioso soprattutto nei dispositivi embedded di piccole dimensioni, come telecamere di sicurezza, droni e dispositivi IoT, ma anche in dispositivi più grandi come le auto su cui ci concentreremo in seguito; i SoC per l'IA facilitano l'esecuzione di inferenze direttamente sul dispositivo, riducendo la necessità di trasferire dati a server remoti, migliorando le prestazioni e la latenza, poiché le operazioni di inferenza sono gestite localmente. Anche nei dispositivi mobili, come smartphone e tablet, l'implementazione di questi SoC consente di eseguire applicazioni di intelligenza artificiale, come il riconoscimento vocale e la visione artificiale, direttamente sul dispositivo. Nei veicoli autonomi, elaborano in tempo reale dati da sensori come telecamere e lidar, contribuendo alla guida autonoma e alla sicurezza stradale [51].

Sebbene questi chip abbiano architetture che differiscono da molti fattori, come l'utilizzo per cui sono stati progettati o il dispositivo su cui sono stati implementati, è possibile andare a delineare una struttura comune formata da:

CPU: La CPU (Unità di Elaborazione Centrale) è il cervello del SoC, responsabile dell'esecuzione delle istruzioni generali del software. Nei SoC per IA, la CPU può essere utilizzata per compiti di controllo del sistema e per operazioni non altamente parallele. Tuttavia, in ambito di intelligenza artificiale, è spesso affiancata da unità specializzate per accelerare specifiche operazioni di machine learning. Solitamente questi processori sono basati su RISC-V, ARM o ISA.

GPU: La GPU (Unità di Elaborazione Grafica) è progettata per gestire operazioni parallele, rendendola ideale per applicazioni di intelligenza artificiale che coinvolgono calcoli matematici intensivi, tipici delle reti neurali [51]. La sua architettura parallela consente di eseguire simultaneamente numerose operazioni, migliorando le prestazioni nei carichi di lavoro di machine learning. Non è raro trovare architetture ottimizzate per l'IA prive di questo componente però, infatti esso può essere sostituito dagli acceleratori hardware .

Acceleratori Hardware: Gli acceleratori hardware, come le NPU, sono componenti specializzate progettate per eseguire specifiche operazioni di machine learning in modo efficiente. Queste unità offrono prestazioni superiori rispetto alle CPU e GPU generali per determinati tipi di carichi di lavoro legati all'IA.

Memoria: La memoria è essenziale per l'archiviazione temporanea dei dati e delle istruzioni. Nei SoC per IA, la gestione della memoria è cruciale per garantire un accesso veloce ai dati utilizzati frequentemente. Essa è solitamente una SRAM [51].

BUS: Il BUS è il sistema di interconnessione che collega tutte le componenti del SoC, consentendo il trasferimento efficiente di dati tra CPU, GPU, acceleratori hardware e memoria. Una progettazione efficace del BUS è fondamentale per evitare congestioni e garantire una comunicazione rapida tra le unità di elaborazione; in questo elemento, la larghezza di banda e la bassa latenza sono essenziali poiché permettono uno scambio di dati maggiore in minor tempo.

Periferiche di I/O: Le periferiche di Input/Output consentono al SoC di comunicare con il mondo esterno, come una SSD o un potenziale processore esterno. Queste includono interfacce per sensori, connessioni di rete, display e altri dispositivi esterni.

La differenza principale tra un SoC progettato specificamente per supportare l'inferenza dell'IA e un SoC generico è spesso la presenza o l'assenza di un acceleratore hardware dedicato alla velocizzazione dei calcoli matriciali necessari per le operazioni di intelligenza artificiale. In questi SoC, come visto in precedenza, spesso si trova un'unità di elaborazione dedicata chiamata NPU (Neural Processing Unit) o un acceleratore hardware specifico per eseguire operazioni matriciali utilizzate nelle reti neurali. Questi acceleratori sono ottimizzati per eseguire rapidamente operazioni di moltiplicazione di matrici e altre operazioni fondamentali per le applicazioni di machine learning. L'inclusione di tali acceleratori consente una maggiore efficienza e velocità nell'esecuzione delle operazioni legate all'IA, migliorando le prestazioni complessive del sistema per carichi di lavoro specifici.

Al giorno d'oggi, su molti SoC, anche se non sono stati esplicitamente sviluppati per supportare l'inferenza dell'IA, sono presenti una varietà di architetture inerenti all'inferenza IA. Un esempio significativo è rappresentato dal Qualcomm Snapdragon 855, che integra un DSP (Digital Signal Processor) chiamato Hexagon 690. Nonostante non sia una NPU dedicata, il DSP svolge un ruolo cruciale nel migliorare le prestazioni in ambito di intelligenza artificiale.

Architetture di un SoC ottimizzato per l'IA e implementato su Smartphone

Apple A17 Pro Bionic

Il cuore tecnologico degli iPhone 15 Pro e iPhone 15 Pro Max è alimentato dal potente Apple A17 Pro, un SoC basato sull'architettura ARM a 64 bit progettato da Apple e prodotto da TSMC. Questo SoC è studiato per offrire prestazioni di livello superiore e avanzate capacità, specialmente nei contesti di intelligenza artificiale.

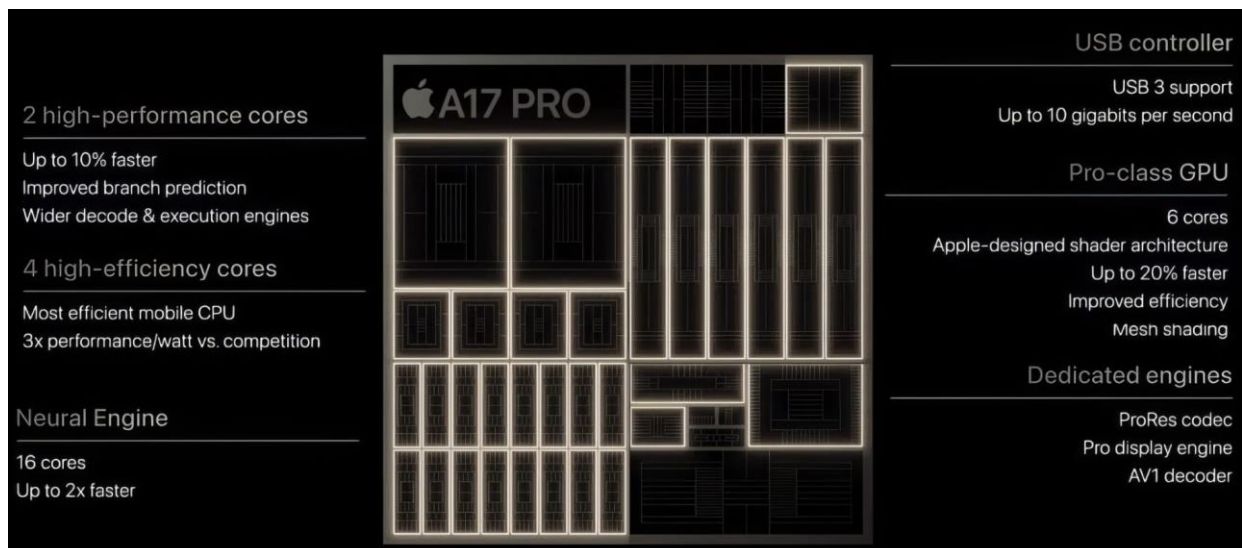


Figura 1: Architettura a blocchi del SoC A17 Pro Bionic [69]

La CPU a sei core, focalizzata sulla massimizzazione delle prestazioni e dell'efficienza energetica, rappresenta una pietra miliare nell'evoluzione dell'architettura Apple. I due core ad alte prestazioni, operanti a una frequenza di 3,78 GHz e i quattro core ad alta efficienza energetica a 2,11 GHz non solo migliorano la velocità [69], offrono un equilibrio ottimale tra potenza computazionale e durata della batteria, elementi fondamentali per un'esperienza d'uso continua e performante nelle applicazioni legate all'IA nei dispositivi mobili.

Grazie alla RAM LPDDR5 da 6 GB o 8 GB è possibile garantire una gestione fluida e veloce di complessi carichi di lavoro e permette una maggiore capacità di multitasking e la gestione efficiente

di dati complessi, essenziale per applicazioni AI-intensive. Per quanto riguarda le memorie in questo SoC sono presenti anche due livelli di cache nella CPU, rispettivamente di 256 KB e 16MB [69].

Ci sono anche 4 percorsi indipendenti nel bus, ciascuno con una capienza di 16 bit e con una larghezza di banda massima pari a 51.2 Gbps.

Nel System on a Chip Apple A17 Pro Bionic, spicca una GPU a 6 core in grado di operare a 2147.2 Gigaflops³ con una frequenza di 1398 MHz. Questo rappresenta un notevole passo avanti, soprattutto nell'ambito dell'inferenza dell'Intelligenza Artificiale su dispositivi mobili. Tuttavia, il componente più notevole per le applicazioni di IA è senza dubbio il Neural Engine, noto anche come ANE (Apple Neural Engine). Questo consiste in un set di 16 core di calcolo altamente specializzati progettati per l'implementazione a basso consumo energetico di reti neurali profonde in grado di raggiungere le 35 TOPS. [31]

Il Neural Engine gioca un ruolo cruciale nell'accelerare in modo significativo gli algoritmi di Machine Learning e Intelligenza Artificiale su dispositivi Apple. Questo componente è responsabile dell'utilizzo di funzionalità come Siri in modalità offline e molte altre applicazioni di IA, inclusa la tecnologia di riconoscimento del volto Face ID.

Il design del Neural Engine è ottimizzato per il calcolo parallelo, dove numerose operazioni, come le moltiplicazioni di matrici eseguite in trilioni di iterazioni, devono avvenire simultaneamente. Inoltre, per accelerare l'inferenza negli algoritmi di IA, il Neural Engine fa uso di modelli predittivi. Dotato della propria cache e supportando solo specifici tipi di dati, il Neural Engine è progettato per massimizzare le prestazioni, contribuendo così all'efficienza complessiva del SoC Apple A17 Pro Bionic. [52]

³ Il termine "Gigaflops" rappresenta la misura della potenza di calcolo di un sistema, indicando miliardi di operazioni in virgola mobile al secondo. Questa unità di misura è comunemente utilizzata per quantificare la capacità di elaborazione dei computer e dei processori, specialmente nei contesti di calcolo scientifico e applicazioni computazionali intensive. La valutazione in Gigaflops è cruciale per valutare le prestazioni di supercomputer, server e dispositivi di calcolo avanzati, riflettendo la loro capacità di eseguire operazioni matematiche complesse in tempi ridotti.

ACCELERATORI HARDWARE

Descrizione degli Acceleratori Hardware per IA

Abbiamo precedentemente parlato dei SoC AI e dei loro vantaggi di utilizzo in applicazioni embedded; in questi chip abbiamo visto che una componente di particolare rilevanza è l'acceleratore hardware, ossia un processore sviluppato appositamente per migliorare le performance di inferenza o apprendimento di una IA.

Questi acceleratori possono trovarsi oltre che all'interno dei SoC anche in Chip dedicati, in modo tale da poter avere dei sistemi più performanti ma sicuramente più ingombranti. I chip vanno distinti in due grandi categorie, quelli utilizzati per l'addestramento e quelli utilizzati per l'inferenza [8] [9].

I chip progettati per la formazione fungono essenzialmente da insegnanti della rete, ed essendo la formazione ad alta intensità di calcolo [14] [9], abbiamo bisogno di chip AI focalizzati sulla formazione progettati per essere in grado di elaborare questi dati in modo rapido ed efficiente. Più potente è il chip, più velocemente la rete apprende. Solitamente questi acceleratori hardware hanno, oltre a una memoria decisamente ampia, delle precisioni numeriche elevate, ad esempio a 32 bit, per evitare di compromettere la precisione del modello, ma non è raro trovare dei dispositivi con precisioni numeriche più basse come a 16 bit, questo perché consente di avere più velocità di calcolo, una memoria più piccola e un risparmio energetico decisamente superiore. In molte applicazioni, soprattutto durante la formazione, è accettabile tollerare una certa quantità di errore numerico. Riducendo la precisione numerica, è possibile ottenere un compromesso tra efficienza computazionale e precisione sufficiente per l'applicazione specifica, ma consente anche di abbattere i costi degli acceleratori. Questi chip hanno inoltre delle memorie decisamente più grandi rispetto ai chip destinati all'inferenza, e questo proprio perché per l'addestramento c'è una grande necessità di dati.

Supponendo ad esempio un IA che si addestra per guidare autonomamente una macchina, essa dovrà avere a disposizione una grandissima quantità di dati come ore e ore di video di macchine che guidano in qualsiasi condizione possibile.

Una volta che una rete è stata addestrata, ha poi bisogno di chip progettati per l'inferenza al fine di utilizzare i dati nel mondo reale, per cose come il riconoscimento facciale, il riconoscimento dei gesti, l'elaborazione del linguaggio naturale, la ricerca di immagini, il filtraggio dello spam ecc..... In questo

caso invece, le architetture degli acceleratori non richiedono grande memoria e precisione come nel caso dell'addestramento [16] [29]; è comune, infatti, trovare acceleratori a 8 bit.

Vale la pena notare, inoltre, che i chip progettati per l'addestramento possono anche eseguire l'inferenza, ma i chip di inferenza non sempre possono eseguire l'addestramento. Queste tipologie di chip possono trovarsi sia su Cloud che su dispositivi Edge, a differenza dei chip per l'addestramento che per motivi di spazio, costo e consumo energetico si trovano per lo più all'interno dei server [14].

Architettura di un Acceleratore ottimizzato per l'addestramento su Cloud

Intel Gaudi2 Habana

Di seguito è riportata l'architettura dell'Intel Gaudi2 Habana, un chip utilizzato per lo più all'interno dei Server.

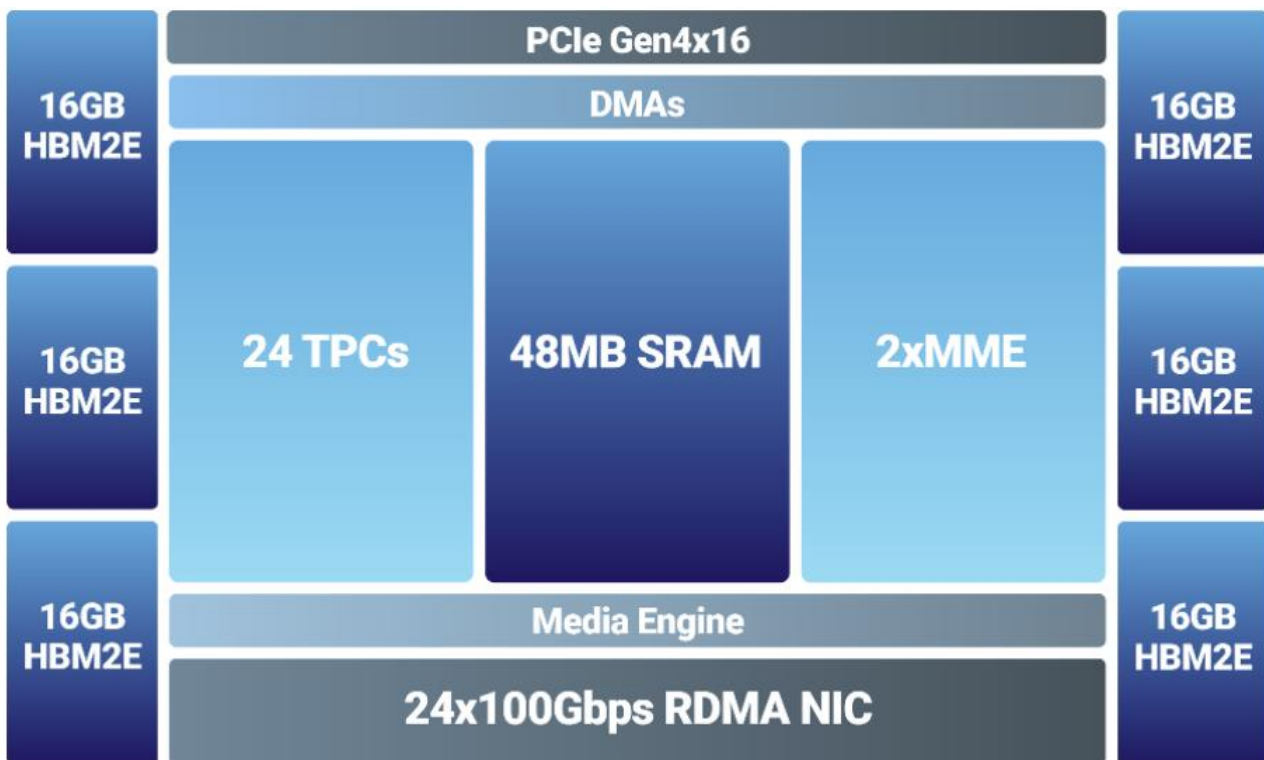


Figura 2: Architettura a blocchi del processore Intel Gaudi2 Habana [70]

Il Chip Gaudi2 Habana di Intel è stato progettato come una microarchitettura per l'accelerazione della formazione delle IA nel data center e presenta una combinazione di componenti specializzati per gestire carichi di lavoro di addestramento di reti neurali nel contesto del deep learning; in particolare presenta 24 TPC e due motori MME. Il TPC (Tensor Processing Cluster) è un processore flessibile e programmabile progettato per gestire in modo efficiente i complessi carichi di lavoro di addestramento delle reti neurali. La sua programmabilità consente agli sviluppatori di adattare gli algoritmi di formazione in base alle esigenze specifiche dell'applicazione, fornendo

flessibilità e ottimizzazione, questo chip infatti come altri mette a disposizione la possibilità di fare operazioni a 32 bit a 16 bit o a 8 bit, in modo da adattare al meglio quelle che sono le mie esigenze con quelli che sono i tempi di sviluppo e le altre risorse a disposizione. Dall'altro lato, il MME (Matrix Math Engine) è un motore dedicato alla moltiplicazione di matrici. Questo componente si occupa di operazioni comuni nel deep learning, come strati completamente connessi, convoluzioni e General Matrix Multiply (GEMM) in batch. La presenza di un motore specializzato per queste operazioni contribuisce a migliorare l'efficienza computazionale complessiva del processore.

Il processore Gaudi2 si distingue per le sue caratteristiche avanzate, offrendo una notevole larghezza di banda di rete, ben 2,4 Terabit e un potente sottosistema di memorie, che comprendendo 96 GB di memorie ad elevata larghezza di banda HBM2E. Inoltre, dispone di 48 MB di SRAM locale con sufficiente larghezza di banda per consentire alle diverse unità, come MME, TPC, DMA e NIC RDMA, di operare in parallelo senza compromettere le prestazioni [71].

Gaudi2 integra nativamente 24 Network Interface Controller (NIC) RoCE V2 RDMA da 100 Gbps direttamente sul chip. Questa integrazione consente una comunicazione efficiente tra più processori Gaudi2, sia attraverso routing diretto che tramite commutazione Ethernet standard. Il risultato è un sistema altamente scalabile e flessibile per gestire i carichi di lavoro di deep learning [30].



Figura 3: Datasheet di Intel Gaudi2 Habana con i miglioramenti rispetto alla versione precedente

[71]

Sicuramente però il fatto di poter collegare più chip di questo genere ad altissima velocità e permettendo così di ottenere un parallelismo impressionante è sicuramente uno dei maggiori punti di forza di Gaudi2.

Di seguito è riportato un metodo di connessione di 8 Gaudi2.

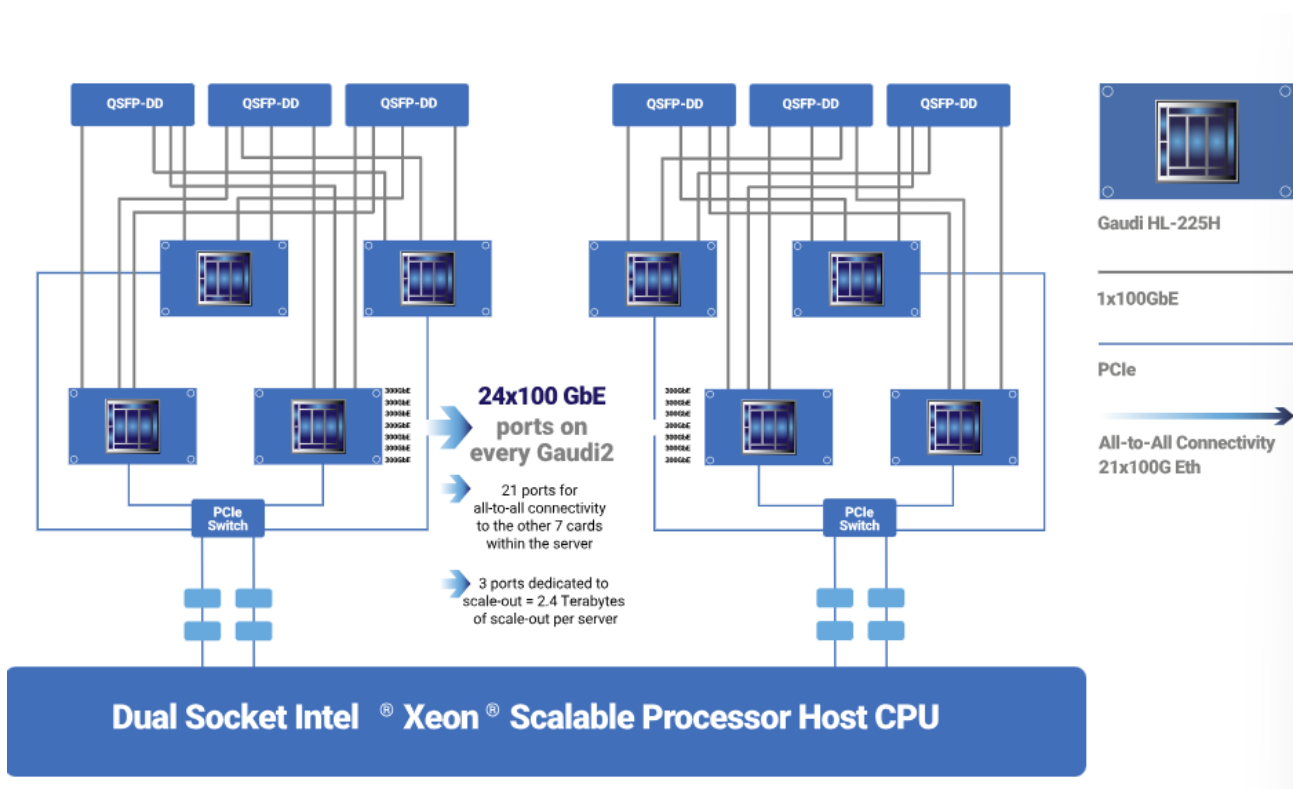


Figura 4: Esempio di collegamento di 8 processori Gaudi2 [72]

Nella figura precedente è possibile osservare come gli otto acceleratori Gaudi2 siano connessi a gruppi di quattro con uno switch PCIe a una CPU Intel Xeon Dual Socket, ossia una famiglia di processori sviluppati per server che offrono un'architettura scalabile, un supporto Dual Socket⁴ ed elevate prestazioni fornite dal grande numero di core, anche diverse centinaia. Questi chip sono connessi tra loro invece tramite le 24 porte RDMA NIC da 100 Gbe ognuna, in particolare 3 sono dedicate allo scale-out, ossia sono destinate alla possibilità di poter espandere ulteriormente questa infrastruttura, e sono collegate a dei QSFP-DD (Quad Small Form-Factor Pluggable Double Density),

⁴ Il termine "Dual Socket" si riferisce a una configurazione hardware di un sistema informatico che include due socket della CPU sulla scheda madre, consentendo l'installazione di due processori distinti. Questa configurazione è spesso utilizzata in server e workstation ad alte prestazioni per aumentare la capacità di elaborazione e gestire carichi di lavoro intensivi. L'utilizzo di due processori consente di distribuire il carico computazionale in modo più efficiente, migliorando le prestazioni complessive del sistema.

ossia uno standard di interfaccia di connettività ottica ad alta velocità, mentre gli altri 21 sono utilizzati per connettere ogni chip con gli altri sette [72].

Di seguito è invece riportata l'architettura del Gaudi 2 MegaPod Architecture, ossia un pod costituito da 8 Gaudi2 Node (quelli visti in precedenza) [72], dove ognuno di questi come abbiamo visto contiene 8 chip Gaudi2. Dalla figura è facile vedere come siano sufficienti 3 switches ARISTA 7060DX4 da 32 porte ognuno per connettere questi server [30] [53].

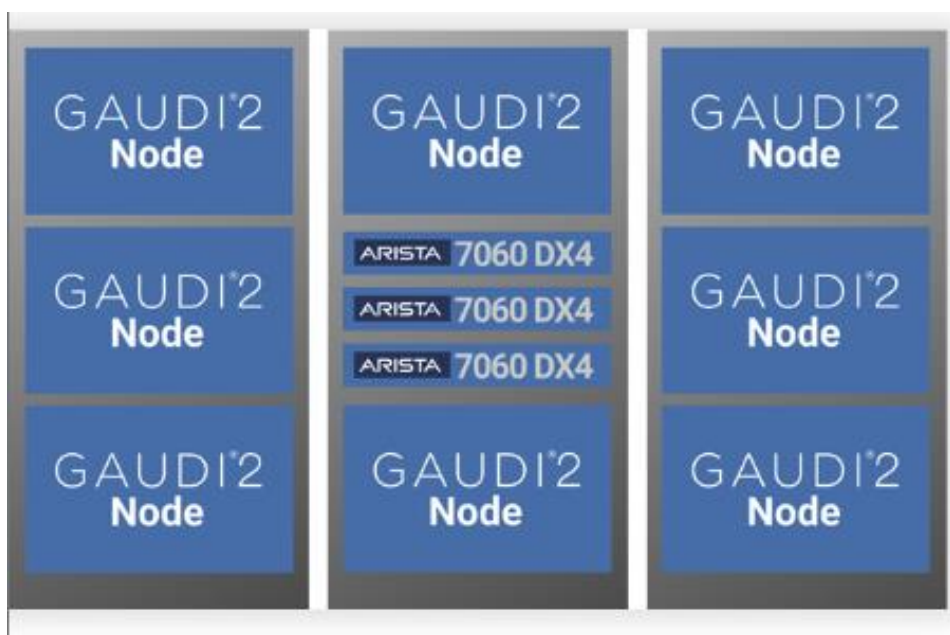


Figura 5: esempio di collegamento di più sistemi Gaudi2 (figura precedente)[72]

Questi acceleratori vengono anche usati all'interno di Amazon AWS in particolare nelle istanze Amazon EC2 DL1, composte da 8 acceleratori Gaudi con a disposizione 4 TB di SSD NVMe e 96 vCPU. Queste istanze forniscono modelli di Deep Learning con costi di addestramento ridotti per i casi d'uso di elaborazione del linguaggio naturale, individuazione di oggetti e riconoscimento di immagini [46].

AWS Trainium

AWS⁵ Trainium rappresenta la seconda generazione di acceleratori di machine learning sviluppati appositamente da AWS per il training di modelli di deep learning caratterizzati da oltre cento miliardi di parametri. Questo acceleratore è stato ottimizzato per supportare l'addestramento di modelli in diverse aree, tra cui l'elaborazione del linguaggio naturale, la visione artificiale e i sistemi di raccomandazione. Le sue applicazioni spaziano su molteplici scenari, inclusi il riepilogo del testo, la generazione di codice, la risposta alle domande, la creazione di immagini e video e il rilevamento di frodi.

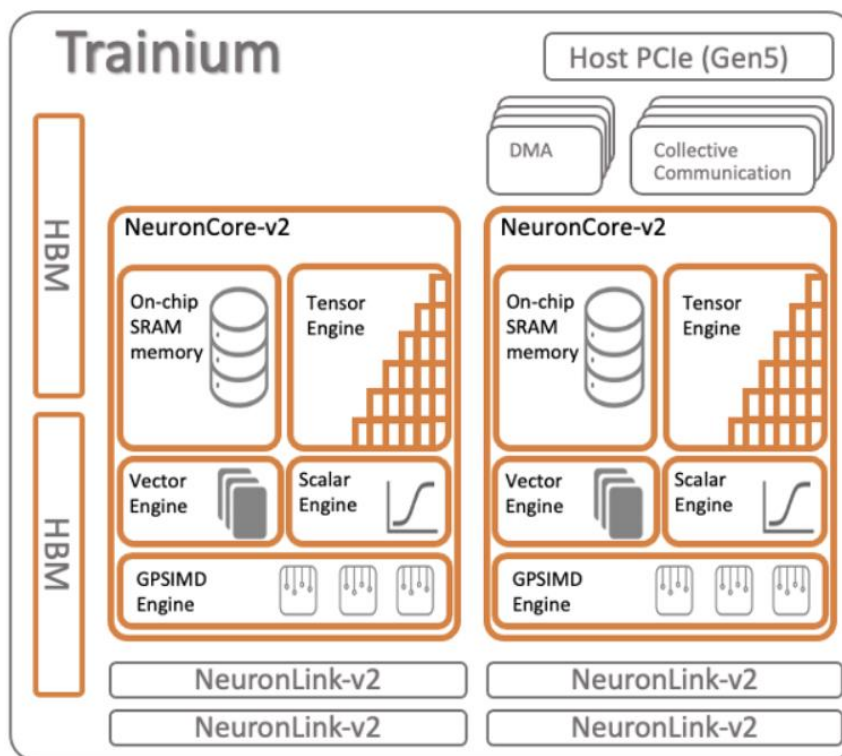


Figura 6: Diagramma a blocchi di AWS Trainium [42]

Ciascun acceleratore Trainium è equipaggiato con due NeuronCore di seconda generazione, progettati appositamente per algoritmi di deep learning, che si distingue per essere un'unità di calcolo eterogenea completamente indipendente [42]. Ogni NeuronCore-v2 include quattro motori principali:

⁵ AWS, acronimo di Amazon Web Services, è una piattaforma di servizi cloud fornita da Amazon e offre una vasta gamma di servizi, tra cui l'archiviazione dei dati, l'elaborazione, la distribuzione di contenuti, l'analisi dati, la gestione delle risorse, e molti altri.

Tensor, Vector, Scalar e GPSIMD, ciascuno con accesso a memoria SRAM gestita tramite software integrato nel chip.

Lo ScalarEngine è specializzato in calcoli scalari, fornendo una notevole potenza di calcolo di 2,9 TFLOPS di calcoli FP32, tre volte più veloce rispetto alla generazione precedente (NeuronCore-v1). Il VectorEngine, ottimizzato per calcoli vettoriali, offre una velocità di 2,3 TFLOPS di calcoli FP32, risultando 10 volte più veloce rispetto a NeuronCore-v1. Il TensorEngine sfrutta un array sistolico, che definiremo tra poco, ottimizzato per i calcoli tensoriali, supportando una vasta gamma di precisioni miste e offrendo oltre 90 TFLOPS di potenza di calcolo FP16/BF16, sei volte più veloce rispetto alla generazione precedente. NeuronCore-v2 introduce anche un nuovo motore denominato GPSIMD-Engine, composto da otto processori generici da 512 bit completamente programmabili. Questi processori possono eseguire codice C lineare e hanno accesso diretto agli altri motori NeuronCore-v2 e alla memoria SRAM incorporata [41].

Per supportare un'efficiente gestione dei dati e il parallelismo dei modelli, ogni acceleratore Trainium è dotato di 32 GB di memoria e una larghezza di banda DMA di 1 TB/sec, supportando la compressione/decompressione in linea della memoria. Questi acceleratori offrono fino a 190 TFLOPS di potenza di calcolo FP16/BF16 e presentano NeuronLink, una tecnologia di interconnessione non bloccante ad altissima velocità che gli permette di poter collegare più chip per ottenere una parallelizzazione delle operazioni decisamente più efficiente [42].

Questi acceleratori sono implementati nelle istanze Trn1 di Amazon Elastic Compute Cloud⁶ (Amazon EC2), ciascuna dotata di un massimo di 16 chip AWS Trainium per garantire elevate prestazioni di elaborazione nei carichi di lavoro di machine learning, un massimo di 128 vCPU e 4 SSD NVMe da 2 TB. Nella figura seguente è possibile vedere come avviene questo collegamento [43].

⁶ Amazon Elastic Compute Cloud (EC2) è un servizio di cloud computing offerto da Amazon Web Services (AWS) che consente agli utenti di eseguire macchine virtuali su richiesta, consentendo agli utenti di aumentare o diminuire le risorse di calcolo in base alle esigenze del carico di lavoro. Gli utenti possono selezionare diverse istanze di macchine virtuali con varie configurazioni di CPU, memoria, archiviazione e rete.

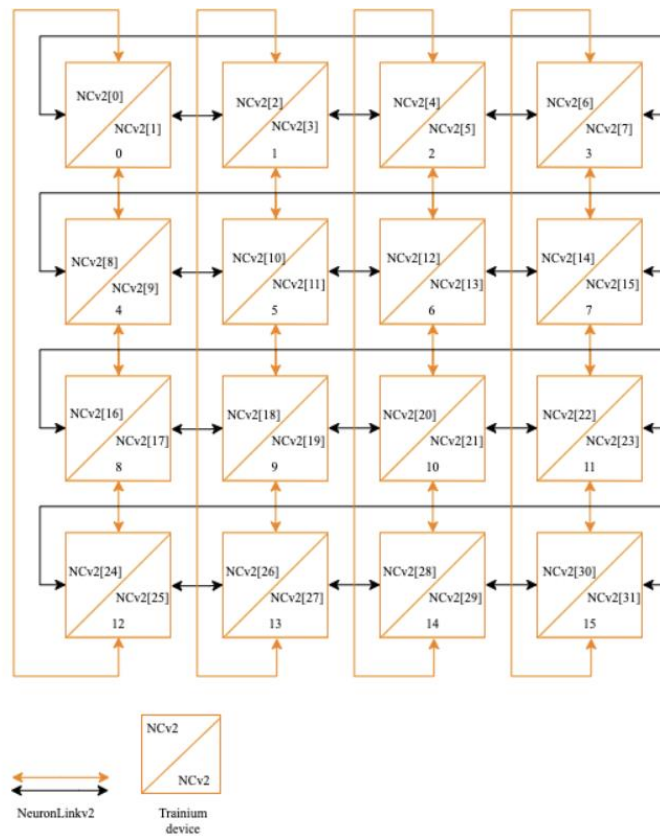


Figura 7: Diagramma a blocchi di Trn1 [43]

Le istanze Trn1/Trn1n sono disponibili anche in un UltraCluster EC2, che consente ai clienti di scalare su oltre 30.000 dispositivi Trainium [40].

Architettura di un Acceleratore ottimizzato per l'inferenza su Cloud

AWS Inferentia 2

Inferentia 2 rappresenta la seconda generazione di acceleratori di inferenza per Cloud sviluppati da AWS che hanno portato a una latenza end-to-end inferiore del 25% e a costi inferiori del 30% rispetto alle istanze basate su GPU. Analogamente a Trainium, è possibile combinare 12 chip di Inferentia 2 per creare un'istanza denominata Inf2.

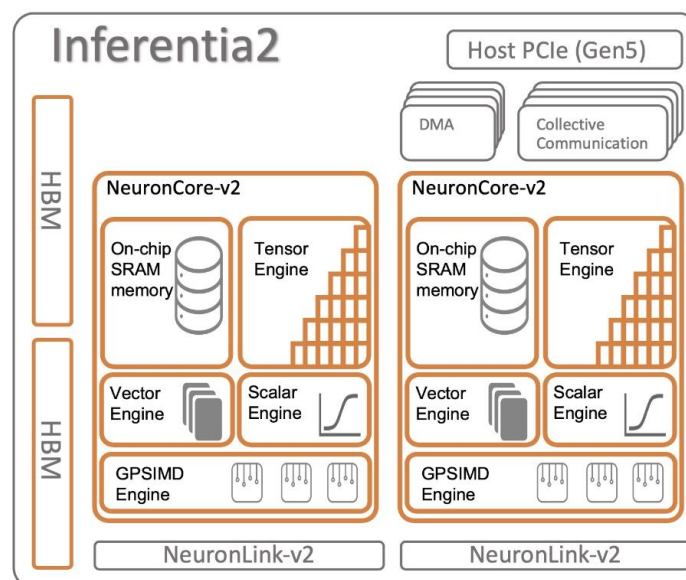


Figura 8: Diagramma a blocchi di AWS Inferentia2 [45]

Ciascun dispositivo Inferentia2 ha un'architettura molto simile a quella dei Trainium; infatti, è costituito da due core NeuronCore-v2, accompagnati da 32 GiB di memoria ad accesso rapido HBM, utilizzata per memorizzare lo stato del modello. Questa memoria è caratterizzata da una notevole larghezza di banda pari a 820 GiB/sec [45].

L'unica differenza con Trainium è che quest'ultimo utilizza ben 4 collegamenti NeuronLink-v2, Inferentia 2 ne utilizza solamente 2, riducendo così il numero di chip in una sola istanza.

Come già detto quindi, collegando tra di loro tramite NeuronLink-v2 12 chip Inferentia2, è possibile ottenere un'istanza chiamata Inf2 con un massimo di 192 vCPU e 768 GiB di memoria. Nell'immagine seguente è possibile vedere come avviene questo collegamento tramite un BUS PCIe [44] [45].

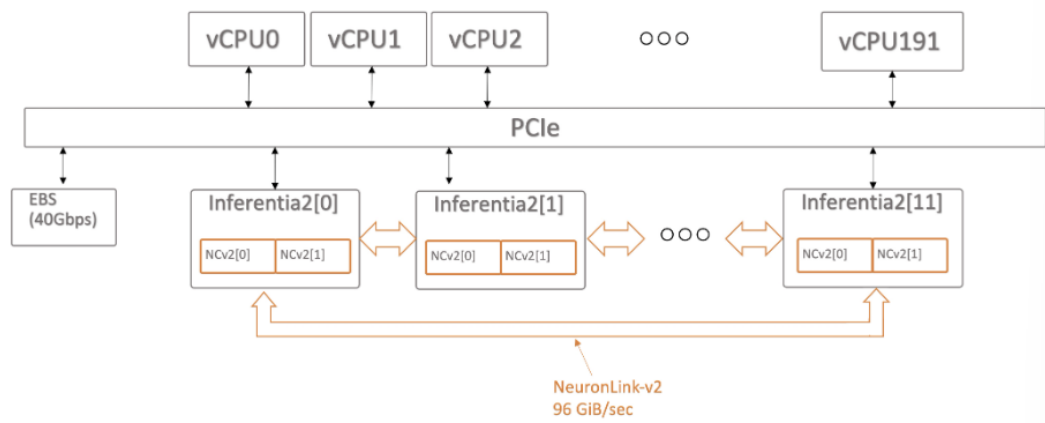


Figura 9: Diagramma a blocchi di Inf2 [44]

Nell'immagine è presente anche Amazon Elastic Block Store (Amazon EBS), ossia un servizio di storage fornito da Amazon Web Services che offre uno spazio di archiviazione persistente per le istanze di Amazon EC2 [44].

Architettura delle GPU

Un' unità di elaborazione grafica (GPU) è un chip di computer che originariamente era responsabile del rendering di immagini, animazioni e video 2D e 3D eseguendo rapidi calcoli matematici, ma ora hanno una gamma di utilizzo più ampia come nel settore delle AI; nel caso in cui la GPU non venga usata per attività grafiche, parleremo di GPGPU (General-Purpose GPU).

Il termine unità di elaborazione grafica è diventato popolare nel 1999 quando Nvidia ha commercializzato la sua GeForce 256 con funzionalità di trasformazione grafica, illuminazione e ritaglio di triangoli. Si tratta di calcoli matematici pesanti, che alla fine aiutano a rendere gli spazi tridimensionali. Questi processori, infatti, eseguono calcoli relativi alla grafica molto rapidamente e in parallelo per consentire un rendering rapido e fluido dei contenuti sullo schermo del computer. Grazie alla GPU, inoltre, la CPU sarà libera di gestire tutto il resto che non è legato all'applicazione grafica.

Le GPU funzionano utilizzando un metodo chiamato elaborazione parallela, in cui divide i problemi complessi in milioni di attività facilitando ulteriormente la ricerca di soluzioni in una sola volta; ed è stata proprio questa elaborazione parallela ad aver permesso un uso massiccio delle GPU in ambito di addestramento e di inferenza di IA.

In particolare, la CPU invia istruzioni alla GPU per disegnare il contenuto grafico sullo schermo. La GPU una volta ricevuta l'istruzione, la divide in migliaia di attività secondarie più piccole ed esegue le istruzioni in parallelo per visualizzare il contenuto sul dispositivo, un processo noto come pipeline grafica o di rendering. La suddivisione in attività più piccole è resa possibile dall'elevato numero di core presenti all'interno di GPU (solitamente migliaia, ma cambiano da GPU a GPU).

Grazie a questo elevato livello di potenza di elaborazione, le GPU sono attualmente una delle scelte più valide sia per il machine learning sia per le attività di inferenza, sia nei dispositivi embedded sia nei Cloud, ma con ovvie differenze architetturali; proprio in questi ultimi, è possibile aumentare la capacità di calcolo aggiungendo più GPU in modo da dividere le attività più grandi in migliaia di attività secondarie più piccole e di elaborarle tutte contemporaneamente.

Tutto questo consente a una GPU di eseguire le operazioni richieste per reti neurali o per gli algoritmi di ML in una frazione di tempo rispetto alle CPU che hanno un numero di core decisamente inferiore, occupandosi di istruzioni più complesse e meno parallelizzabili. I modelli di DL o ML, infatti, possono essere addestrati più velocemente semplicemente eseguendo tutte le operazioni

contemporaneamente anziché attendere che vengano completate sequenzialmente una dopo l'altra come nelle CPU.

Un'altra caratteristica che le rende ottimali per le implementazioni di IA è la loro larghezza di banda della memoria, essa infatti è indispensabile per alimentare simultaneamente i suoi core e spostare grandi quantità di dati velocemente con il processore e la memoria.

Anche per quanto riguarda l'inferenza, l'uso di una GPU migliora di tantissimo l'elaborazione dell'output richiesto; ad esempio, Nvidia dichiara che l'uso di una GPU A100 80GB è in grado di migliorare le performance dell'IA di ben 249 volte rispetto al solo utilizzo di una sola CPU. [27]

Le GPU, pur vantando notevoli vantaggi, presentano alcuni svantaggi che meritano considerazione; come il consumo energetico elevato presente specialmente nelle GPU di fascia alta, richiedono una quantità significativa di energia, il che è uno svantaggio soprattutto per i data center, i quali hanno al loro interno un gran numero di GPU.

Ad esempio, Perlmutter, il più grande supercomputer di IA del 2021, aveva a disposizione 6159 GPU A100, una quantità esorbitante che genera ingenti consumi e altrettanto calore; anche il raffreddamento, dunque, diventa un aspetto cruciale soprattutto per i Data Center, aumentando la complessità delle architetture [24].

Una GPU è costituita da diversi componenti chiave che lavorano insieme. Ogni componente svolge un ruolo fondamentale nel garantire un'elaborazione efficiente e prestazioni ottimali.

Shader Core

Il cuore vitale di una GPU è costituito dagli Shader Core, noti anche come processori stream o core CUDA nelle GPU Nvidia. Queste piccole unità di elaborazione svolgono il ruolo cruciale di eseguire istruzioni e compiere calcoli, operando in parallelo per gestire grandi quantità di dati contemporaneamente. La presenza di un maggior numero di shader core in una GPU aumenta la sua capacità di gestire calcoli complessi [18].

Uno shader core stesso contiene diversi elementi. Contiene decodificatori, dispatcher, raccoglitori di operandi, raccoglitori di risultati e altro ancora. Ma l'elemento più importante, è sicuramente l'Unità Aritmetico Logica (ALU) [28]. Le ALU sono gli elementi costitutivi più fondamentali di una GPU e sono l'unità di base che esegue effettivamente le operazioni matematiche richieste come parte di un programma shader. Nelle GPU Nvidia ci sono 2 ALU su ogni CUDA core: una ALU a virgola mobile FP32 e una ALU intera.

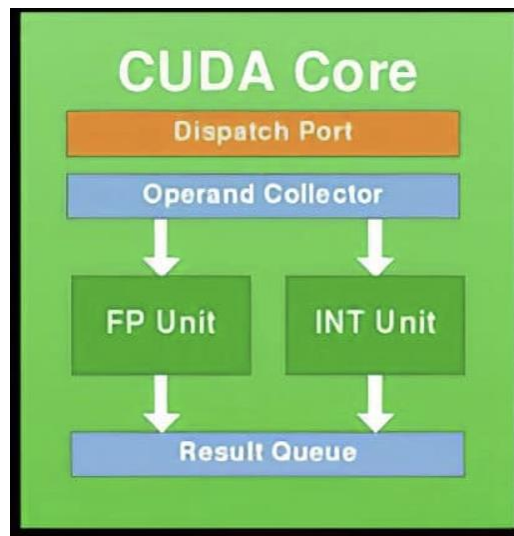


Figura 10: Schema a blocchi di un CUDA core [54]

Gli shader core sono organizzati in insiemi chiamati Streaming Multiprocessor (SM), ma il tessuto connettivo di una GPU si estende oltre gli SM poiché più di questi collaborano per formare un Cluster [19].

Per quanto riguarda le applicazioni di IA, tuttavia, gli shader core sono più veloci dei core CPU comuni quando si tratta di elaborare i numeri, ma non sono ancora la soluzione ideale. Questo perché non sono mai stati concepiti per essere utilizzati in quel modo, ma sono stati realizzati appositamente per l'elaborazione grafica. Proprio per questo, quando le GPU hanno iniziato a essere utilizzate per carichi di lavoro di intelligenza artificiale e machine learning, Nvidia ha introdotto i core Tensor, utilizzati anche da Google nelle sue TPU. Mentre i core CUDA possono eseguire solo un'operazione per ciclo di clock, i core Tensor possono gestire più operazioni, offrendo loro un incredibile incremento delle prestazioni. Fondamentalmente, tutto ciò che i Tensor Core fanno è aumentare la velocità di moltiplicazione della matrice. Questo aumento della velocità di calcolo va a scapito della precisione, poiché i core CUDA sono significativamente più accurati. Detto questo, quando si tratta di addestrare modelli di machine learning, i core Tensor sono molto più efficaci in termini di velocità di calcolo e costo complessivo; quindi, la perdita di accuratezza viene spesso trascurata [54].

TMU (Texture Mapping Unit)

Le TMU (Texture Mapping Units) svolgono un ruolo essenziale nell'aggiunta di dettagli realistici alle superfici degli oggetti 3D. Queste unità sono incaricate di prelevare i dati delle texture dalla memoria e di applicarli in modo accurato ai triangoli o ai poligoni pertinenti. In questo modo, le TMU contribuiscono a creare ombreggiature e dettagli che conferiscono una qualità grafica superiore. Le GPU possono integrare più TMU per gestire in modo efficiente le complesse texture necessarie per ottenere risultati visivi di alta qualità nell'ambito della grafica avanzata. Tuttavia, va precisato che questa componente non è di particolare rilievo per le applicazioni di IA.

ROP (Raster Operation Unit)

Le ROP sono responsabili di convertire i dati prodotti dagli shader in pixel visualizzati sullo schermo. Queste unità traducono in dettaglio i risultati generati dagli shader, determinando il colore, la profondità e altri attributi di ciascun pixel nell'immagine finale. Sono responsabili della scrittura dei dati pixel nella memoria, ma come il componente precedente non è decisamente importante per l'efficacia di una GPU nell'ambito dell'IA.

Frame Buffer

Il frame buffer rappresenta un'area dedicata della memoria video destinata a conservare l'immagine conclusiva che sarà proiettata sullo schermo. All'interno di questa memoria, sono archiviati i valori dei pixel, le informazioni cromatiche, i dati sulla profondità e altri attributi relativi a ciascun pixel della visualizzazione. Questo spazio di memoria viene continuamente aggiornato durante l'elaborazione dei dati da parte della GPU, riflettendo così le modifiche in tempo reale e costituendo la base per l'output visivo finale [18].

Gerarchia di Memoria

I registri, che rappresentano la forma più veloce di memoria, sono utilizzati per archiviare temporaneamente dati e registri di istruzioni per le unità di elaborazione. Un altro tipo di memoria, sono le cache di primo livello (L1), che, situata direttamente sul chip della GPU, fornisce una memoria veloce e a bassa latenza per le unità di elaborazione. Ogni core ha la sua cache L1, quindi la Cache L1 è spesso dedicata e non condivisa tra i vari core.

Proseguendo nella gerarchia, si incontra la cache di secondo livello (L2), con una capacità di memorizzazione più ampia rispetto alla L1. La sua funzione principale è ridurre la latenza di accesso ai dati rispetto alla memoria principale, contribuendo ad aumentare l'efficienza complessiva della GPU. La Cache L2 viene utilizzata per memorizzare dati e istruzioni condivise tra i core all'interno della stessa unità (come SM o cluster) [19].

Inoltre, va detto che a disposizione della GPU vi è una memoria globale (VRAM) che rappresenta la sua principale area di archiviazione, destinata a immagini, texture, frame buffer e altri dati cruciali per il rendering grafico. Pur avendo una capacità maggiore rispetto alle cache, l'accesso alla VRAM è più lento [19].

Infine, in alcune situazioni specifiche, le GPU possono attingere alla memoria di sistema del computer (RAM) per archiviare dati aggiuntivi, quando la VRAM risulta insufficiente o per scopi particolari come l'addestramento di modelli di machine learning, in questo caso parleremo di memoria unificata, che permette anche a più GPU di collaborare, come nei Data Center.

La memoria, in particolare quella principale, riveste un ruolo cruciale nel determinare le capacità di una GPU nell'ambito dell'intelligenza artificiale. Essa deve essere in grado di sostenere elevati carichi di dati e operazioni parallele necessarie per l'addestramento di reti neurali e altri compiti di machine learning. La gestione efficiente della memoria è fondamentale per accelerare la velocità di accesso ai dati durante le elaborazioni, contribuendo così all'ottimizzazione delle prestazioni complessive della GPU nell'esecuzione di compiti legati all'IA. Una memoria di alta larghezza di banda è particolarmente vantaggiosa in questo contesto, garantendo che la GPU possa elaborare rapidamente grandi quantità di dati.

GPU NVIDIA H100 Tensor Core

La GPU Nvidia H100 Tensor Core, basata sull'architettura Hopper, rappresenta la nona generazione di unità di elaborazione grafica per data center sviluppata da Nvidia. Questa mostruosa GPU, con ben 14592 CUDA core, è progettata per offrire prestazioni notevolmente superiori, fino a 30 volte rispetto alla generazione precedente A100 Tensor Core, particolarmente ottimizzata per carichi di lavoro di intelligenza artificiale e calcolo ad alte prestazioni su larga scala [33].

Uno dei principali impieghi dell'H100 è nel campo dell'addestramento e dell'esecuzione di modelli di intelligenza artificiale, come evidenziato dall'utilizzo da parte di OpenAI per alimentare il modello di linguaggio ChatGPT nel supercomputer Azure. Aziende come Meta hanno adottato H100 per la creazione di supercomputer avanzati, come nel caso di Grand Teton.

Form Factor	H100 SXM	H100 PCIe
FP64	34 teraFLOPS	26 teraFLOPS
FP64 Tensor Core	67 teraFLOPS	51 teraFLOPS
FP32	67 teraFLOPS	51 teraFLOPS
TF32 Tensor Core	989 teraFLOPS*	756teraFLOPS*
BFLOAT16 Tensor Core	1,979 teraFLOPS*	1,513 teraFLOPS*
FP16 Tensor Core	1,979 teraFLOPS*	1,513 teraFLOPS*
FP8 Tensor Core	3,958 teraFLOPS*	3,026 teraFLOPS*
INT8 Tensor Core	3,958 TOPS*	3,026 TOPS*
GPU memory	80GB	80GB
GPU memory bandwidth	3.35TB/s	2TB/s
Decoders	7 NVDEC 7 JPEG	7 NVDEC 7 JPEG
Max thermal design power (TDP)	Up to 700W (configurable)	300-350W (configurable)
Multi-Instance GPUs	Up to 7 MIGS @ 10GB each	
Form factor	SXM	PCIe Dual-slot air-cooled
Interconnect	NVLink: 900GB/s PCIe Gen5: 128GB/s	NVLINK: 600GB/s PCIe Gen5: 128GB/s
Server options	NVIDIA HGX™ H100 Partner and NVIDIA-Certified Systems™ with 4 or 8 GPUs NVIDIA DGX™ H100 with 8 GPUs	Partner and NVIDIA-Certified Systems with 1-8 GPUs
NVIDIA AI Enterprise	Add-on	Included

Figura 11: Datasheet GPU NVIDIA H100 Tensor Core [25]

Le caratteristiche chiave che contribuiscono alla potenza di H100 includono i tensor core di quarta generazione, con prestazioni fino a sei volte superiori rispetto al modello A100. L'aggiunta di una nuova unità TMA (tensor memory accelerator) consente un trasferimento efficiente di grandi blocchi di dati tra la memoria globale e quella condivisa, ottimizzando ulteriormente le prestazioni.

L'architettura di memoria è potenziata dalla tecnologia HBM3⁷, ossia una VRAM che offre una notevole larghezza di banda di 3 TB/sec. La cache L2 da 50 MB contribuisce a ottimizzare gli accessi ripetuti a grandi porzioni di modelli e set di dati, riducendo i tempi di accesso al sottosistema di memoria HBM3 [25].

La connettività multi-GPU è abilitata dalla tecnologia Nvidia NVLink di quarta generazione e la nuova tecnologia NVSwitch di terza generazione, offrendo un aumento significativo nella larghezza di banda e accelerando le connessioni tra più GPU, infatti, i modelli di intelligenza artificiale da trilioni di parametri per attività richiedono mesi per essere addestrati, anche sui supercomputer. Per accorciare i tempi di sviluppo alla velocità aziendale e completare la formazione in poche ore è necessaria una comunicazione continua ad alta velocità tra tutte le GPU in un cluster di server.

Per affrontare carichi di lavoro di grandi dimensioni, i nuovi NVLink e NVSwitch sono stati progettati per abilitare l'aggregazione di 8 GPU H100, formando l'HGX H100 8-GPU insieme a 4 NVSwitch. In questa configurazione, ogni GPU è dotata di multiple porte NVLink e si connette a tutti e quattro gli switch. Ciascuno di questi è uno switch completamente non bloccante che collega tutte e otto le GPU H100 Tensor Core. Questa topologia completamente connessa consente a ciascuna GPU H100 di comunicare simultaneamente con tutte le altre, operando a una velocità bidirezionale di NVLink di 900 gigabyte al secondo (GB/s)[33] [34]. Questo valore rappresenta oltre 7 volte la larghezza di banda dell'attuale bus PCIe Gen5.

⁷ HBM, acronimo di High Bandwidth Memory, è una tecnologia di memoria avanzata utilizzata in alcuni dispositivi elettronici, in particolare nelle schede grafiche e nei processori accelerati. Al contrario delle memorie tradizionali come la GDDR, l'HBM impiega uno stack di chip di memoria verticalmente all'interno di un package, riducendo la latenza e migliorando la larghezza di banda dei dati tra la memoria e il processore. L'adozione di HBM è stata particolarmente significativa nelle schede grafiche ad alte prestazioni, consentendo un accesso più rapido ai dati e un miglioramento delle prestazioni complessive del sistema.

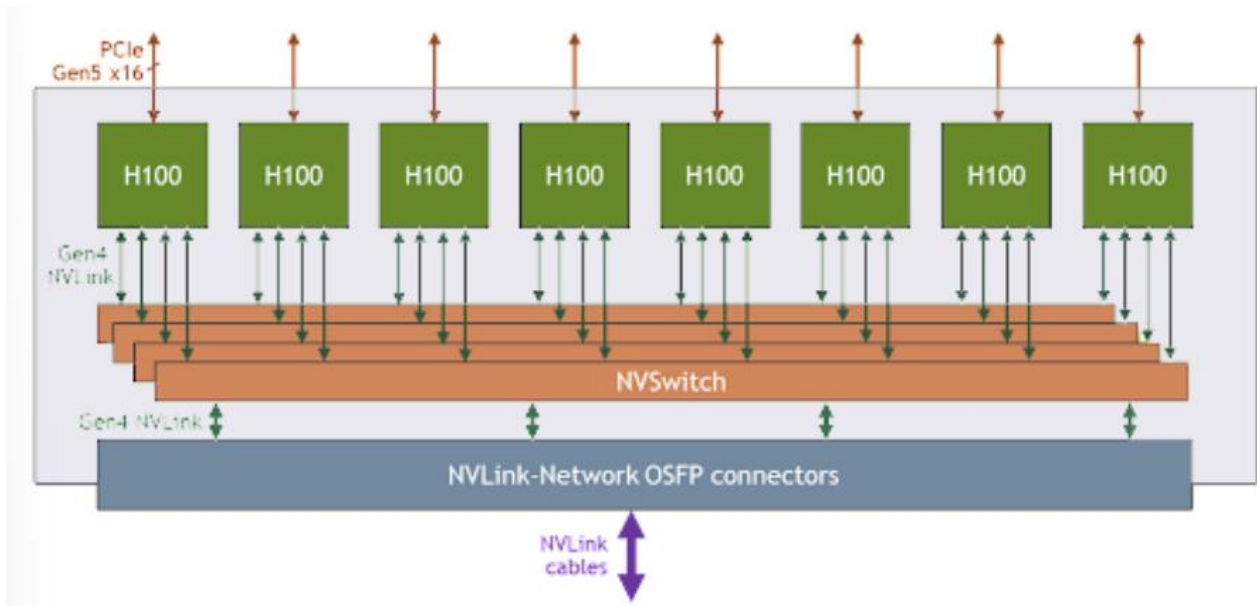


Figura 12: Diagramma a blocchi di alto livello dell'HGX H100 8-GPU con supporto di rete NVLink [34]

L'evoluzione della classe di applicazioni emergenti, come l'HPC exascale e i modelli di intelligenza artificiale con trilioni di parametri, richiede tempi di addestramento notevolmente ridotti, passando da mesi a poche ore. Questo richiede una comunicazione continua ad alta velocità tra tutte le GPU in un cluster di server. Gli HGX H100 8-GPU con NVLink supportano un dominio più ampio grazie alla nuova componente NVLink-Network, con cui è possibile collegare fino a un massimo di 256 domini GPU, facilitando la creazione di cluster estremamente potenti per gestire complessi carichi di lavoro di intelligenza artificiale e HPC su larga scala [33] [34].

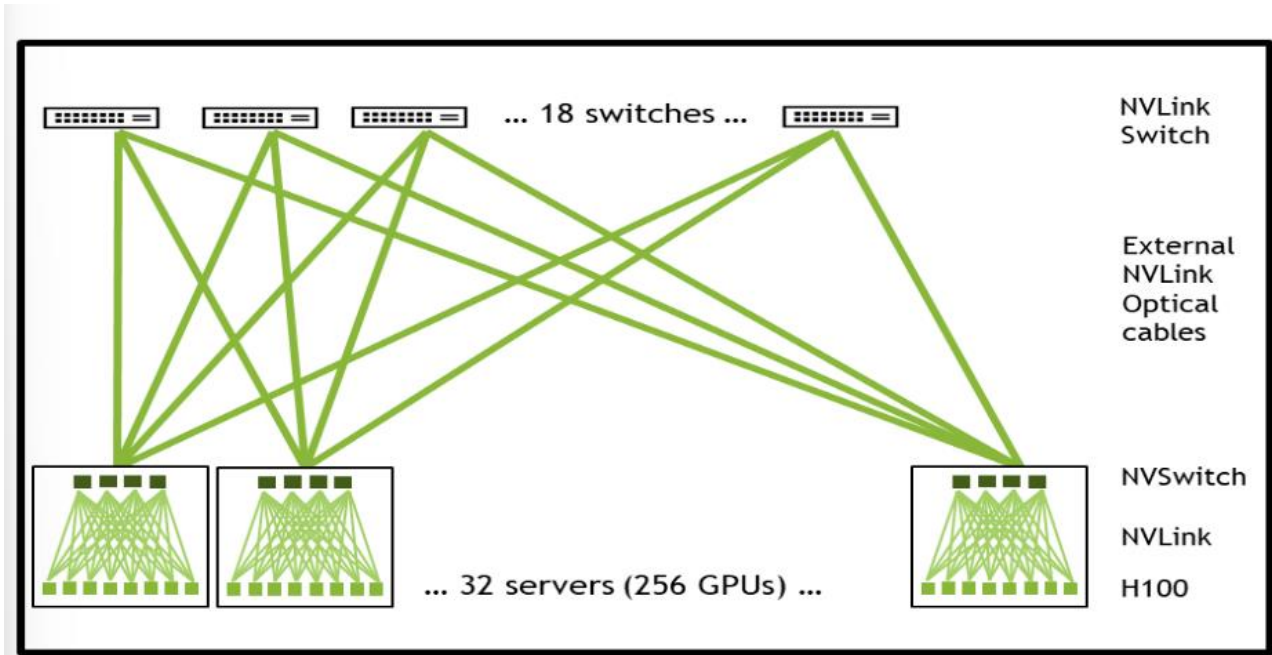


Figura 13: Pod GPU 256 H100 [34]

Attraverso test di prestazioni, è stato dimostrato che un sistema con 256 GPU H100 può raggiungere notevoli livelli di potenza computazionale. I risultati includono 1024 petaflop con FP8, indicando la capacità del sistema di eseguire circa 1024 quadrilioni di operazioni in virgola mobile al secondo con una precisione a 8 bit. Inoltre, il sistema può generare 512 petaflop con FP16, rappresentando la capacità di eseguire circa 512 quadrilioni di operazioni in virgola mobile al secondo con una precisione a 16 bit. Infine, sono stati registrati 15 petaflop con FP64, indicando la capacità del sistema di eseguire circa 15 quadrilioni di operazioni in virgola mobile al secondo con una precisione a 64 bit, adatta per applicazioni che richiedono una maggiore precisione numerica [34].

GPU NVIDIA TITAN RTX

La GPU NVIDIA TITAN RTX è un processore che seppur appartiene alla categoria TITAN e quindi è a tutti gli effetti una GPU da gaming, è stata progettata anche per supportare la ricerca sull'intelligenza artificiale su dispositivi embedded, a differenza delle H100; e proprio per questo è possibile notare moltissime differenze con la GPU studiata precedentemente, come il numero di CUDA core, Tensor core o la grandezza e la larghezza di banda della memoria, al fine di permettere a questa GPU di essere implementata in dispositivi che hanno capacità differenti.

SPECIFICATIONS	
GPU Memory	24 GB GDDR6
Memory Interface	384-bit
Memory Bandwidth	Up to 672 GB/s
NVIDIA CUDA® Cores	4,608
NVIDIA Tensor Cores	576
NVIDIA RT Cores	72
Single-Precision Performance	16.3 TFLOPS
Tensor Performance	130 TFLOPS
NVIDIA NVLink	Connects 2 TITAN RTX GPUs
NVIDIA NVLink Bandwidth	100 GB/s (bidirectional)
System Interface	PCI Express 3.0 x 16
Power Consumption	280 W
Thermal Solution	Active
Form Factor	4.4" H x 10.5" L, Dual slot, full height
Display Connectors	3x DisplayPort, 1x HDMI, 1x USB Type-C
Max Simultaneous Displays	4x 4096 x 2160 @ 120 Hz, 4x 5120 x 2880 @ 60 Hz, 2x 7680 x 4320 @ 60 Hz
Encode/Decode Engines	1x encode, 1x decode
VR Ready	Yes
Graphics APIs	Microsoft DirectX 12 API³, Vulkan API⁴, OpenGL 4.6⁴
Compute APIs	CUDA, DirectCompute, OpenCL™

Figura 14: Scheda tecnica delle GPU NVIDIA TITAN RTX [73]

Costruito sull'architettura Turing, è dotato di 4608 CUDA cores, 576 Tensor Core di precisione mista che forniscono fino a 130 teraFLOPS (TFLOPS) per la formazione e l'inferenza sull'apprendimento profondo e 72 core RT per accelerare il ray tracing. Accanto alla GPU ci sono 24 GB GDDR6 di

memoria per l'addestramento di reti neurali con batch di grandi dimensioni, a riguardo c'è da sottolineare anche la larghezza di banda della memoria che può raggiungere 672 GB/s, garantendo un accesso veloce ed efficace ai dati [32] [73].

Esattamente come nelle GPU H100 poi, è possibile grazie alla tecnologia NVLink connettere due GPU per ottenere performance migliori [55].

Architettura di un'istanza che utilizza GPU per l'addestramento su Cloud

Nei capitoli precedenti abbiamo visto come Amazon metta a disposizione nei propri Cloud delle istanze per l'apprendimento delle IA come Trn1 o DL1e istanze per l'inferenza come Inf2 che utilizzano dei particolari chip costruiti appositamente per questo; tuttavia, Amazon ha creato anche altre istanze come P4d che utilizzano delle GPU A100 Tensor Core.

Amazon EC2 P4d

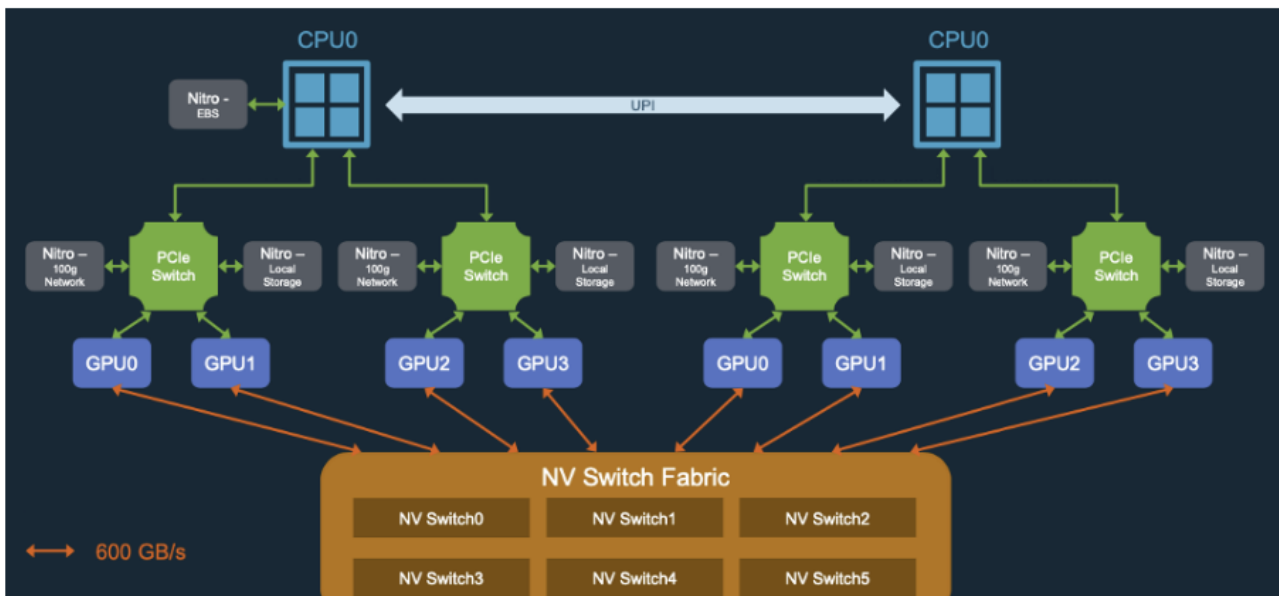


Figura 15: Diagramma a blocchi di un'istanza P4d di AWS [48]

Questa risorsa è dotata di processori Intel Cascade Lake 8275CL dual socket per un totale di 96 vCPU a 3,0 GHz con 1,1 TB di RAM e 8 TB di memoria locale NVMe con una velocità di lettura superiore ai 16 GB/s, ma il cuore pulsante dei P4d sono le 8 GPU NVIDIA Tesla A100 collegate tramite NVSwitch [47] [48], dove ogni GPU integra 56 miliardi di transistor ed è equipaggiata con 80GB di memoria HBM2e, il chip prevede 108 SM (Streaming Multiprocessor) per un totale di 6192 Cuda Core e 432 Tensor Core; gli 80GB di memoria HBM2e sfruttano un BUS a 6144 bit che, sempre nella variante PCIe, garantiscono una banda passante di picco che tocca i 2 TB/s [26] [20].

**NVIDIA A100 TENSOR CORE GPU SPECIFICATIONS
(SXM4 AND PCIE FORM FACTORS)**

	A100 80GB PCIe	A100 80GB SXM
FP64	9.7 TFLOPS	
FP64 Tensor Core	19.5 TFLOPS	
FP32	19.5 TFLOPS	
Tensor Float 32 (TF32)	156 TFLOPS 312 TFLOPS*	
BFLOAT16 Tensor Core	312 TFLOPS 624 TFLOPS*	
FP16 Tensor Core	312 TFLOPS 624 TFLOPS*	
INT8 Tensor Core	624 TOPS 1248 TOPS*	
GPU Memory	80GB HBM2e	80GB HBM2e
GPU Memory Bandwidth	1,935GB/s	2,039GB/s
Max Thermal Design Power (TDP)	300W	400W***
Multi-Instance GPU	Up to 7 MIGs @ 10GB	Up to 7 MIGs @ 10GB
Form Factor	PCIe dual-slot air cooled or single-slot liquid cooled	SXM
Interconnect	NVIDIA® NVLink® Bridge for 2 GPUs: 600GB/s ** PCIe Gen4: 64GB/s	NVLink: 600GB/s PCIe Gen4: 64GB/s
Server Options	Partner and NVIDIA- Certified Systems™ with 1-8 GPUs	NVIDIA HGX™ A100- Partner and NVIDIA- Certified Systems with 4,8, or 16 GPUs NVIDIA DGX™ A100 with 8 GPUs

Figura 16: Scheda tecnica GPU Tensor Core A100 [74]

NPU

Una Neural Processing Unit (NPU) rappresenta la famiglia di tutti quei microprocessori progettati per ottimizzare le prestazioni delle reti neurali artificiali e degli algoritmi di apprendimento automatico [21]. Come abbiamo visto le applicazioni di Intelligenza artificiale che siano di inferenza o di formazione, possono essere accelerate tramite una parallelizzazione massiccia di processori come le GPU. Tuttavia, gli avanzamenti tecnologici recenti hanno sollevato la necessità di un'accelerazione con un focus particolare sull'efficienza energetica, specialmente in risposta alle crescenti applicazioni nei dispositivi di elaborazione mobile. Mentre le GPU offrono il vantaggio del calcolo parallelo, spesso richiedono la collaborazione di una CPU. La costruzione di modelli di reti neurali e la gestione dei flussi di dati rimangono compiti che la CPU svolge ancora. La GPU, inoltre, può presentare sfide in termini di consumo energetico e ingombro fisico, diventando meno adatta per dispositivi compatti e mobili. Di conseguenza, le NPU con chip dedicati sono diventate cruciali, offrendo dimensioni ridotte, basso consumo energetico, elevate prestazioni computazionali e una notevole efficienza complessiva.

La NPU opera emulando il funzionamento di neuroni e sinapsi umani a livello di circuito, elaborando direttamente su larga scala attraverso un set di istruzioni di deep learning. Queste istruzioni consentono di completare l'elaborazione di un insieme di neuroni con un'efficienza superiore rispetto a CPU e GPU. Le NPU integrano memoria e calcolo attraverso pesi sinaptici, ottimizzando così l'efficienza operativa e a differenza delle migliaia di istruzioni richieste dai processori CPU e GPU per elaborare i neuroni, le NPU possono completare tali operazioni con poche istruzioni, fornendo chiari vantaggi in termini di efficienza di elaborazione nel contesto del deep learning. Risultati sperimentali indicano che le NPU superano le GPU in termini di prestazioni, con un consumo energetico comparabile [22].

Sul mercato sono disponibili diverse varianti di NPU, come ad esempio la Tensor Processing Unit (TPU) di Google che analizzeremo nel capitolo successivo o le Apple NPU viste nei capitoli precedenti [21].

Architettura delle TPU

La Tensor Processing Unit (TPU) rappresenta un acceleratore hardware specializzato progettato da Google per ottimizzare le operazioni di machine learning [22]. Queste unità hanno trovato impiego in diverse applicazioni, principalmente all'interno dell'ecosistema di Google, integrandosi in servizi come Google Foto per migliorare l'elaborazione delle immagini e RankBrain per ottimizzare gli algoritmi di ricerca o anche per addestrare la chatbot di Google Bard.

Le TPU nascono come una tecnologia da utilizzare all'interno dei sistemi cloud, avendo anche delle dimensioni difficilmente trasferibili in sistemi periferici, ma Google in risposta alla crescente enfasi sull'elaborazione locale dei dati al di fuori dei data center centralizzati, ha introdotto nel 2018 l'Edge TPU. Questa variante è progettata per l'Edge computing come suggerisce il nome, caratterizzata da dimensioni ridotte e un consumo energetico inferiore. L'Edge TPU è orientato verso applicazioni di machine learning direttamente sui dispositivi, consentendo calcoli senza la necessità di una connessione cloud costante.

Sono due le motivazioni fondamentali per cui le TPU hanno delle prestazioni maggiori rispetto all'accelerazione CPU/GPU: quantizzazione e matrice sistolica.

La quantizzazione è il primo passo dell'ottimizzazione, che utilizza numeri interi a 8 bit per approssimare numeri in virgola mobile a 16 o 32 bit. Ciò può ridurre la capacità di memoria e le risorse di calcolo richieste. L'array sistolico contribuisce in modo determinante all'efficienza del TPU grazie alla sua naturale compatibilità con la manipolazione della matrice unita al fatto che il calcolo nelle reti neurali può essere rappresentato come operazioni di matrice.

Da una prospettiva globale, il nucleo dell'intero TPU è la Matrix Multiply Unit, che è un 256×256 array sistolico composto da più celle di calcolo. Ciascuna cella riceve un parametro di peso insieme a un segnale di ingresso alla volta ed esegue l'accumulo dei propri prodotti. Una volta che tutti i pesi e i segnali di input sono stati propagati alle celle vicine, inizia immediatamente il ciclo successivo di iterazione. Con questo schema di calcolo, l'intera moltiplicazione della matrice viene effettivamente completata dalla collaborazione di tutte le celle di calcolo. L'array sistolico di MXU contiene $256 \times 256 = 65.536$ ALU, il che significa che il TPU può elaborare 65.536 moltiplicazioni e addizioni di numeri interi a 8 bit per ciclo o meno se si utilizzano approssimazioni a 16 bit [35]. A causa dell'architettura sistolica, i dati di input vengono effettivamente riutilizzati più volte. Pertanto, può raggiungere un throughput più elevato consumando meno larghezza di banda della memoria [38].

Questa caratteristica consente alle TPU di elaborare grandi quantità di dati in parallelo, accelerando notevolmente le attività di machine learning.

Le TPU incorporano anche la memoria a larghezza di banda elevata (HBM), che assicura un accesso più rapido ai dati rispetto alle architetture di memoria tradizionali. Questa scelta di progettazione garantisce un'alimentazione coerente dei dati alla MXU, riducendo i potenziali colli di bottiglia e massimizzando la velocità effettiva computazionale [35].

Tra le caratteristiche delle TPU figura anche l'aritmetica a precisione ridotta, che consente di gestire calcoli con precisione inferiore senza compromettere l'accuratezza dei modelli di apprendimento automatico. Questo approccio mira a migliorare ulteriormente la velocità e l'efficienza complessiva delle TPU.

D'altra parte, le unità di elaborazione grafica sono state inizialmente concepite per la manipolazione della grafica del computer e presentano una struttura parallela adatta per algoritmi che lavorano su blocchi di dati di grandi dimensioni, tipici dei carichi di lavoro di intelligenza artificiale. Google ha progettato il TPU come un processore a matrice, invece che come un processore generico in modo che fosse specializzato nei carichi di lavoro della rete neurale. Ciò risolve il problema di Google relativo all'accesso alla memoria che rallenta GPU e CPU, costringendole a utilizzare più potenza di elaborazione.

Sebbene le GPU siano versatili e in grado di gestire una vasta gamma di attività, potrebbero non sempre eguagliare le prestazioni grezze delle TPU nelle specifiche operazioni di machine learning. Tuttavia, per compiti più ampi e determinate architetture di reti neurali, le GPU rimangono la scelta preferita.

Le TPU sono caratterizzate da una memoria a larghezza di banda elevata, garantendo un accesso rapido ai dati durante i calcoli. Tuttavia, la quantità di memoria su chip potrebbe essere limitata, rappresentando una restrizione per modelli con ingenti requisiti di memoria. D'altro canto, le GPU generalmente dispongono di un pool di memoria più ampio, rendendole adatte per attività che richiedono un ampio accesso o archiviazione dei dati [10].

La versione V4 delle Tensor Processing Unit (TPU) presenta una configurazione avanzata composta da 4 chip. Ogni chip è dotato di 2 Tensor Core e dispone complessivamente di 32 GB di memoria HBM2, caratterizzata da un'ampia larghezza di banda di memoria pari a 1200 GBps [39].

All'interno di ciascun core, troviamo unità di elaborazione vettoriale (VPU), unità scalari e quattro MXU. È importante notare che tali operazioni utilizzano la precisione BF16⁸, che impiega più bit per rappresentare l'esponente rispetto ai tradizionali floating-point a 16 bit. Questa scelta di precisione consente una moltiplicazione interna efficiente, mentre l'accumulo avviene con una precisione di tipo FP32 [36] [37].

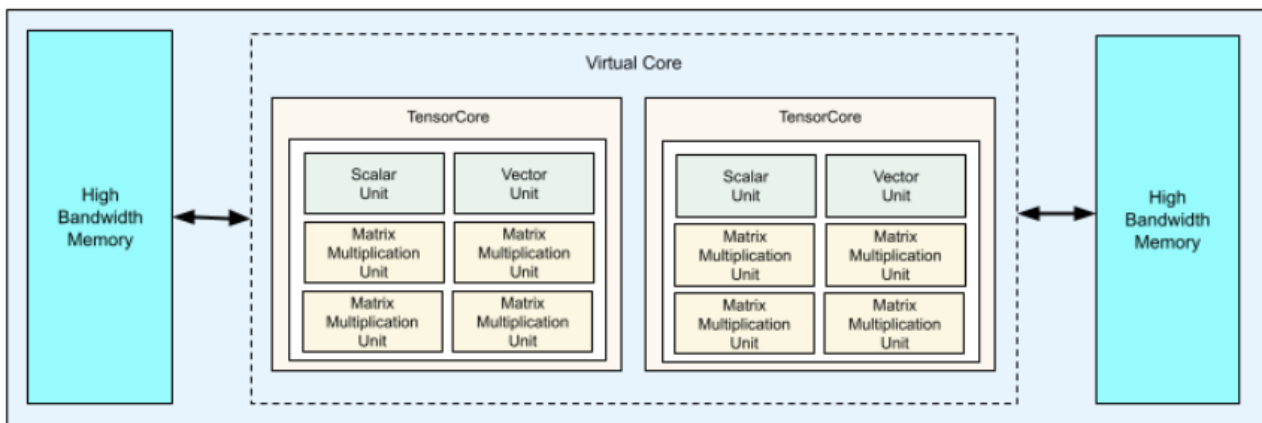


Figura 17: Schema a blocchi di un TPU v4 [39]

È possibile organizzare più dispositivi TPU in configurazioni ad alta velocità, creando così un pod TPU. Ad esempio, un pod TPU V3 può ospitare fino a 256 dispositivi, con un totale di 2048 core TPU V3 e 32 terabyte di memoria. Nel caso della TPU V4, la tecnologia avanzata di interconnessione Google TPU Network consente di formare pod con un numero ancora maggiore di chip, arrivando a 4096. I supercomputer TPU V4 sono resi disponibili per ricercatori e sviluppatori di intelligenza artificiale presso il cluster ML di Google Cloud situato in Oklahoma [39].

Il processo di elaborazione dei dati avviene attraverso un flusso ben definito. L'host TPU invia i dati attraverso una coda di feed, dalla quale la TPU carica i dati nella memoria HBM. Una volta

⁸ BF16 (BFloat16) e FP16 (Half Precision) sono formati di dati a virgola mobile utilizzati nei calcoli numerici, soprattutto nel contesto degli acceleratori di intelligenza artificiale. BF16 è un formato a virgola mobile a 16 bit progettato per offrire un compromesso tra la precisione numerica e l'efficienza di memorizzazione. È particolarmente utile nelle applicazioni di deep learning e machine learning. FP16, o Half Precision, è un altro formato a virgola mobile a 16 bit che rappresenta i numeri con una precisione ridotta rispetto ai formati a 32 bit più comuni. Anche FP16 è ampiamente utilizzato nelle applicazioni di intelligenza artificiale, in quanto offre un miglior compromesso tra prestazioni e requisiti di memoria rispetto ai formati a virgola mobile più precisi.

completato il calcolo, i risultati vengono inseriti nella coda di outfeed. L'host TPU recupera quindi i risultati dalla coda di outfeed e li archivia nella sua memoria.

Per eseguire operazioni sulla matrice, la TPU carica i parametri dalla memoria HBM nell'unità di moltiplicazione matriciale (MXU). Successivamente, la TPU carica i dati dalla memoria HBM. Durante ciascuna operazione di moltiplicazione, i risultati vengono passati all'accumulatore della successiva operazione di moltiplicazione. L'output finale è la somma di tutti i risultati delle moltiplicazioni tra dati e parametri [36]. Ciò avviene senza la necessità di accessi continui alla memoria durante il processo di moltiplicazione della matrice, contribuendo a ottimizzare le prestazioni complessive del sistema.

CONFRONTO TECNICO DELLE ARCHITETTURE OTTIMIZZATE PER L'IA

In questo capitolo ho voluto confrontare le architetture viste nei capitoli precedenti per evidenziare quelle che sono le caratteristiche distintive di questi chip. Di seguito è riportata una tabella in cui per ogni chip vengono indicate: la tipologia e la dimensione della memoria usata, la larghezza di banda della memoria, il numero di operazioni al secondo, il prezzo, il consumo energetico, il numero di dispositivi che si possono collegare per aumentare le prestazioni e il tipo di utilizzo che ne viene fatto.

Tabella di confronto di più architetture							
	Memory	Memory Bandwidth	TOPS	Price	Energy Consumption	Horizontal Scalability	Usage
Apple A17 Pro Bionic	LPDDR5 da 8GB	51.2 Gbps	35 TOPS (INT32)	N/A	11W	1	Iphone 15 - 15Pro
Intel Gaudi2 Habana	6 HBM2e da 16GB	2.4 Tbps	400 TFLOPS (FP16)	7.500,00 €	600W	256	Cloud o Server Privati
AWS Trainium	2 HBM da 16GB	1 Tbps	190 TFLOPS (FP16)	N/A	N/A	16	AWS Cloud
AWS Inferentia 2	2 HBM da 16GB	820 Gbps	190 TFLOPS (FP16)	N/A	N/A	12	AWS Cloud
GPU NVIDIA H100 SXM	HBM3 da 80GB	3.35 Tbps	2000 TFLOPS (FP16)	42.000,00 €	700W	256	Cloud o Server Privati
GPU NVIDIA TITAN RTX	GDDR6 da 24 GB	672 Gbps	130 TFLOPS (FP16)	2.500,00 €	280W	2	Dispositivi Embedded
GPU NVIDIA A100 SXM	HBM2e da 80 GB	2.04 Tbps	624 TFLOPS (FP16)	16.500,00 €	400W	16	Cloud o Server Privati
Google TPUv4	HBM2 da 32 GB	1.2 Tbps	275 TFLOPS (FP16)	N/A	192W	4096	Google Cloud

Figura 18: Tabella riassuntiva dei chip AI

Cominciando dalla memoria, è interessante notare come tutti i dispositivi che vengono utilizzati sui server cloud fanno uso di memorie HBM, mentre le TITAN RTX e l' A17 Pro usano rispettivamente le GDDR e le LPDDR; questo è sicuramente dovuto al fatto che le HBM come già anticipato precedentemente offrono delle prestazioni sicuramente maggiori delle altre due alternative, ma a un costo maggiore e a un consumo energetico decisamente più elevato che non consentirebbe a queste memorie di essere utilizzate in dispositivi piccoli come gli smartphone o dispositivi embedded come dei PC. Un altro fattore che incide su questa scelta è poi quello delle applicazioni IA che i vari dispositivi richiedono, poiché un semplice smartphone viene utilizzato al massimo per l'inferenza di semplici IA non computazionalmente onerose come quelle che vengono usate sui server cloud e quindi tutta la larghezza di banda che le HBM offrono non sarebbe neanche sfruttata al massimo.

Collegato al tipo di memoria e alla loro dimensione è la larghezza di banda delle memorie che nei dispositivi per cloud server è dell'ordine dei Tera per secondo, mentre per i chip per dispositivi embedded è dell'ordine dei Giga per secondo.

La colonna dei TOPS è una colonna indispensabile per capire la pura potenza dei vari chip ; in questa colonna emerge ancora di più la potenza delle GPU NVIDIA H100 SXM che con le loro 2000 tops rappresentano l'opzione più prestazionale riuscendo a compiere in un singolo secondo più di tre volte le operazioni delle GPU NVIDIA A100. Questo dato si scontra però con il prezzo di queste GPU pari a ben 42000€.

Nella colonna del prezzo va precisato tuttavia che non è stato possibile trovare dati ufficiali sul prezzo delle soluzioni AWS e Google in quanto sono utilizzabili solo all'interno dei loro Server e non sono acquistabili separatamente; anche l'Apple A17 Pro Bionic non è acquistabile separatamente poiché viene venduto all'interno degli smartphone Apple.

Le H100 però non primeggiano solo nella colonna del numero di operazioni al secondo, ma anche in quella del consumo energetico che si assesta secondo la sua scheda tecnica intorno ai 700 watt.

Infine, un'ultima precisazione va fatta sulla colonna della scalabilità dei chip, infatti sarà possibile andare a scalare su più dispositivi all'interno dei cluster, per esempio le AWS Trainium possono scalare su oltre 1000 dispositivi, ma nella tabella è indicato il numero 16 perché come visto nei capitoli precedenti, AWS mette a disposizione istanze composte da 16 chip.

ESEMPIO DI UN'ARCHITETTURA PER L'IMPLEMENTAZIONE DI UN'AUTO A GUIDA AUTONOMA

In questo capitolo ho deciso di affrontare in maniera critica quella che potrebbe essere l'architettura necessaria per implementare un'intelligenza artificiale in grado di poter guidare autonomamente un veicolo.

Il primo passo da affrontare è il livello di guida che vogliamo far sviluppare alla nostra autovettura; difatti esistono vari livelli definiti dalla Society of Automotive Engineers (SAE)⁹ International, dallo 0 al 5, dove il livello zero rappresenta la guida manuale e quindi nessuna automazione, mentre il quinto rappresenta una guida autonoma totale, ossia non esiste il bisogno di un conducente umano a bordo [61].

Considerando questa scala, credo sia interessante affrontare una vettura che soddisfi proprio il livello 5; è importante sottolineare che, al momento, nessuna vettura è effettivamente in grado di soddisfare completamente questo livello di automazione su strade pubbliche. Tuttavia, il Livello 5 rappresenta l'apice delle aspirazioni dell'industria automobilistica e simboleggia un futuro in cui la guida quotidiana si trasformerebbe radicalmente.

Il Livello 5 di automazione implica che un veicolo sia completamente autonomo, senza la necessità di interventi umani in qualsiasi circostanza o condizione stradale. Attualmente, le vetture autonome in circolazione sono spesso limitate a scenari specifici o aree geografiche predefinite, e l'aspirazione al Livello 5 è una visione che sfida le concezioni tradizionali della guida.

Raggiungere il Livello 5 significherebbe superare le sfide più complesse associate alla guida autonoma, come la gestione di situazioni di traffico estremamente complesse, condizioni meteorologiche avverse e interazioni con una vasta gamma di utenti della strada, tra cui pedoni e ciclisti. Questa forma di automazione totale comporterebbe una rivoluzione nei concetti di mobilità e

⁹ La Society of Automotive Engineers (SAE) International è un'organizzazione globale senza scopo di lucro che si focalizza sulla promozione dell'innovazione e dello sviluppo nell'industria automobilistica e dell'ingegneria della mobilità. Gli standard SAE sono ampiamente accettati e utilizzati nell'industria automobilistica e forniscono una base comune per la progettazione, la produzione e il test di veicoli e componenti correlati.

potrebbe portare a un cambiamento fondamentale nel modo in cui le persone concepiscono e vivono gli spostamenti quotidiani.

Nel delineare l'architettura hardware per lo sviluppo di un sistema di guida autonoma di alto livello, è imperativo considerare attentamente la divisione tra Edge Computing e Cloud Computing. Tale suddivisione si riflette in due dispositivi distinti, ognuno con un ruolo specifico nell'ecosistema dell'intelligenza artificiale.

Per l'addestramento dell'IA, un approccio orientato al Cloud Computing risulta essenziale, come già visto nei capitoli precedenti. La complessità e le dimensioni dei modelli richiedono notevoli risorse di storage e capacità computazionali, facilmente accessibili attraverso infrastrutture cloud. In questa fase cruciale, la disponibilità di grandi quantità di dati e la possibilità di distribuire l'elaborazione su server potenti favorisce un apprendimento efficace e approfondito.

D'altro canto, per l'inferenza, la fase in cui l'IA prende decisioni in tempo reale durante la guida autonoma, il focus si sposta verso l' Edge Computing. La necessità di ridurre al minimo la latenza diventa fondamentale per garantire risposte immediate e sicure durante le operazioni di guida. L'inferenza sul dispositivo Edge, direttamente installato sull'automobile, evita i rischi associati alla trasmissione dati a lunga distanza, come latenze elevate, potenziali errori di trasmissione o mancanza di connessione.

STORAGE PER ADDESTRAMENTO DI MACCHINA A GUIDA AUTONOMA

Analisi sul miglior dispositivo di storage per l'addestramento di questa
specifica IA

Iniziando dall'apprendimento quindi sarà necessario avere in primo luogo un grande set di dati costituiti per lo più da video di guida, immagini di segnali, ma anche informazioni provenienti da radar, giroscopi, ecc..... in modo tale da poter addestrare l'IA nella maniera più completa possibile, prevedendo tutti i possibili scenari. In genere per avere una buona base di inizio per questi progetti così ambiziosi saranno necessari migliaia di terabyte di dati.

Una volta accumulati questi dati, tuttavia, sarà necessario portare dei continui aggiornamenti, sia perché magari c'è la necessità di studiare una nuova dinamica, sia perché può esserci la necessità di voler completamente cambiare un comportamento inizialmente dato per corretto, e questo può essere fatto sia nelle fasi di addestramento preliminare e quindi quando ancora il software non è stato reso disponibile, ma anche quando l'IA comincerà a funzionare sui dispositivi Edge. È ovvio che l'aggiornamento del set di dati non è una pratica complessa; tuttavia, potrebbe capitare che all'interno dei dati utilizzati per l'addestramento manchino delle casistiche ed ecco perché è importante riportare un metodo fornito da Tesla per tenere aggiornati i dati di addestramento; in ogni auto Tesla, infatti, è equipaggiato un computer Full Self-Driving (FSD) che gestisce simultaneamente due sistemi FSD distinti. Un computer FSD guida attivamente il veicolo quando è attivo il pilota automatico, mentre l'altro opera in una "modalità nascosta" [59].

La "modalità nascosta" funziona in modo simulato come se stesse effettivamente controllando l'auto. Quando il conducente adotta un comportamento diverso da quello previsto o quando la rete neurale segnala incertezze in determinati scenari, tali eventi vengono registrati come imprecisioni. Queste imprecisioni vengono archiviate nella memoria del sistema, consentendo a Tesla di raccogliere dati retroattivamente.

Nel caso in cui Tesla identifichi delle imprecisioni simili in circostanze analoghe, può esaminare situazioni di guida analoghe che si sono verificate su altre vetture della flotta Tesla, anche se l'errore specifico non è stato rilevato. La società può quindi compilare esempi contestuali simili. Sfruttando

questi nuovi set di dati accuratamente etichettati, Tesla può successivamente riformare la sua rete neurale per gestire meglio le situazioni in cui si sono verificate le imprecisioni [59].

Per far sì quindi di avere uno spazio di archiviazione sufficiente, che ovviamente ha la necessità di variare nel tempo, la scelta più ovvia è quella di un servizio Cloud come ad esempio EC2 di AWS, che offre una scalabilità praticamente illimitata, con i dati accessibili da qualunque luogo e al contempo mantenuti in sicurezza a costi contenuti per le grandi big tech.

Con la scelta di un sistema di archiviazione su Cloud, anche la scelta delle dimensioni necessarie per contenere i dati diventa marginale; questo perché sarà possibile espandere o ridurre lo spazio di archiviazione in qualunque momento a seconda delle necessità e questo fa sì che anche i costi cambino di conseguenza.

RISORSE COMPUTAZIONALI PER ADDESTRAMENTO DI UNA IA PER UNA MACCHINA A GUIDA AUTONOMA

Analisi del miglior dispositivo hardware per l'addestramento di questa
specifica IA

Una volta scelta la destinazione dei dati, sempre per quanto riguarda l'addestramento, sarà opportuno scegliere delle risorse computazionali all'avanguardia, in grado di poter svolgere moltissimi calcoli e addestrare velocemente l'IA. In generale, progetti di questo genere richiedono diversi mesi, se non anni, di lavoro.

Prima di andare a determinare quali siano le architetture migliori per l'addestramento di una IA di questo genere, è necessario andare a capire con quale formato andare a fare questo addestramento; ossia se addestrare l'IA con 32 bit, 16 bit o 8 bit. La differenza sta nel fatto che meno bit di precisione vengono usati, meno l'IA sarà precisa nelle sue scelte avendo un margine d'errore più elevato, ma di contro i costi saranno decisamente più bassi, poiché sarà necessario meno spazio di archiviazione e anche i calcoli saranno svolti più velocemente.

Considerando però il nostro caso di studio e sottolineando il fatto che le automobili dovranno essere in grado di scegliere l'alternativa più giusta in qualsiasi momento, poiché ne viene la vita delle persone, mi sembra normale optare per l'opzione di addestramento a 32 bit o 16 bit .

Prima di andare a fare uno studio sulle migliori architetture per addestramento attualmente in commercio, e quindi fare un confronto tra GPU, TPU e GAUDI2 (3 delle opzioni studiate in precedenza per l'addestramento di IA su Cloud) credo sia opportuno andare a delineare quella che è la miglior GPU per questa casistica, in modo tale da fare un confronto più dettagliato.

Considerando quindi che molto probabilmente l'addestramento avverrà, come già detto, su Cloud il confronto verrà fatto unicamente sulle capacità computazionali delle due GPU NVIDIA attualmente reputate più potenti; la A100 e l'H100 SXM.

Per quanto riguarda le prestazioni, la GPU NVIDIA H100 SXM ha raggiunto una velocità di 2,2x per Bf16 rispetto alle A100; tuttavia, una scoperta interessante emerge quando si confrontano i costi di gestione di queste GPU nel cloud. CoreWeave fissa il prezzo delle GPU H100 SXM a 4,76

dollari/ora/GPU, mentre il prezzo della A100 SXM da 80 GB ottiene 2,21 dollari/ora/GPU. Sebbene l'H100 sia 2,2 volte più costoso, le prestazioni compensano, con conseguente minor tempo per addestrare un modello e un prezzo inferiore per il processo di addestramento. Ciò rende intrinsecamente H100 più attraente per i ricercatori e le aziende che desiderano addestrare grandi modelli di Intelligenza Artificiale e rende più fattibile la scelta della GPU H100, nonostante l'aumento dei costi.

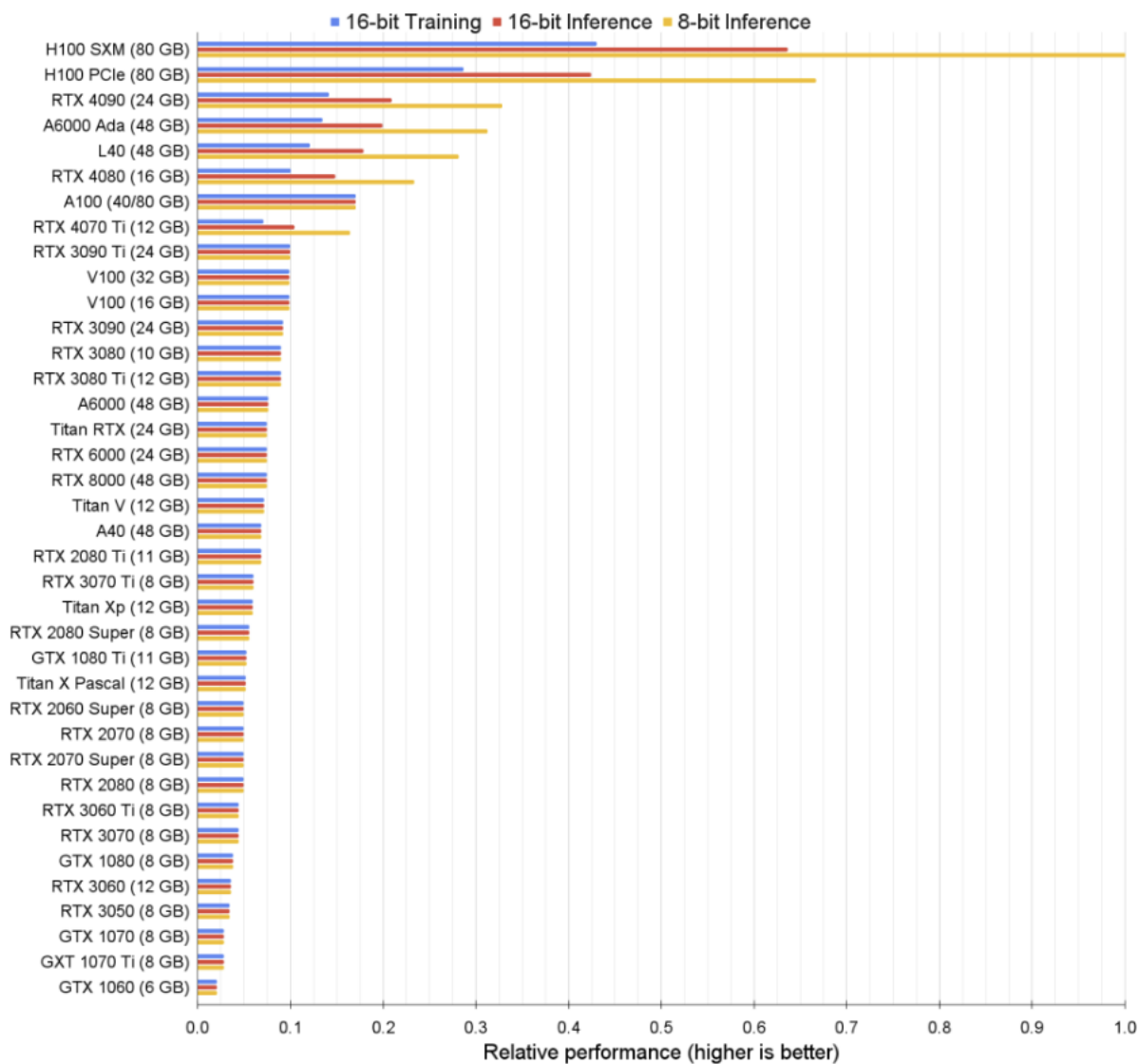


Figura 19: Prestazioni delle migliori GPU sul mercato [75]

Considerata ora la H100 SXM la miglior GPU per l'addestramento su cloud per le nostre esigenze e andando anche a fare un confronto con le altre due alternative, ossia le TPU v5 di Google e il Gaudi2 di Intel, emerge che le H100 sono più veloci delle altre due opzioni sia per l'addestramento di IA per il riconoscimento di immagini (StableDiffusion), sia per IA per il rilevamento di oggetti (ResNeSt), sia per IA per la trascrizione del linguaggio umano (BERT). Di seguito sono riportati i risultati del benchmark di MLCommons (Figura 19).

Intel® Gaudi®2 performance advances strengthen competitive price-performance vs. H100

- Gaudi2 performance on ResNet near that of H100.
- H100 with FP8 outperformed Gaudi2 with BF16 on BERT.
- Vs. TPU, Gaudi2 delivered 3x performance on GPT-3.
- Given its significantly lower server cost vs. H100 server cost, Intel Gaudi2 delivers price-performance advantage vs. H100 across models.

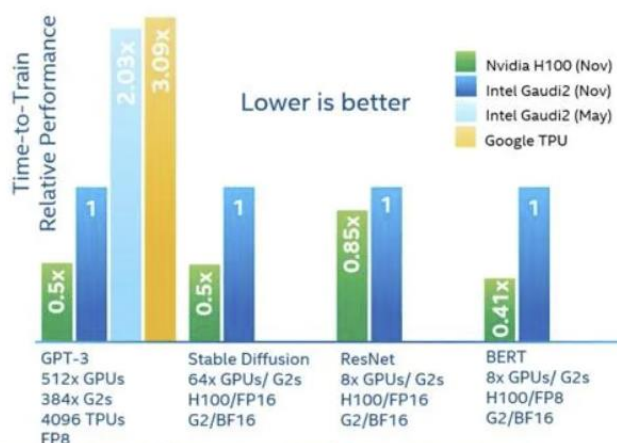


Figura 20: Benchmark di NVIDIA H100, Google TPU e Intel Gaudi2 su diverse applicazioni AI [60]

Ovvio è poi che la scelta di un'architettura piuttosto che un'altra è dettata anche dalle offerte messe a disposizione dalle aziende cloud, a quel punto saranno i vari analisti delle aziende che mettendo sul piatto tutte queste offerte e facendo una comparazione costi/ricavi considerando anche il fattore tempo saranno in grado di scegliere l'architettura migliore; basandoci tuttavia esclusivamente sulle prestazioni del singolo componente, considerando anche il benchmark riportato, la scelta migliore ricade sulle H100 SXM [60].

ARCHITETTURA DELLA MACCHINA A GUIDA AUTONOMA

Analisi dell'architettura di un'autovettura di livello 5

Come prima osservazione per quando riguarda questo tipo di architettura c'è da determinare quelle che potrebbero essere tutte le periferiche indispensabili come Radar, Lidar, Telecamere, GPS, ADAS¹⁰. Grazie a tutte queste componenti e a tutti quei componenti e ad altri componenti che permettano alla macchina di sterzare in maniera autonoma, frenare, accelerare, ecc.... il computer di bordo insieme all'IA presviluppata su Cloud potrà essere in grado di pilotare autonomamente un'autovettura.

Tutti questi sensori e le unità di controllo necessarie, anche dette Zone ECU¹¹, dovranno essere collegate al computer di bordo, in modo tale da poter elaborare gli input e determinare gli output; collegamento che può essere fatto ad esempio tramite Automotive Ethernet, ossia una tecnologia di comunicazione basata sulla rete Ethernet, progettata specificamente per soddisfare le esigenze di connettività avanzata nei veicoli moderni [62]. Questa tecnologia è diventata sempre più popolare nell'industria automobilistica grazie alla sua capacità di supportare una vasta gamma di applicazioni e servizi in un ambiente automobilistico sempre più connesso e complesso. È attualmente usata anche per i sistemi ADAS anche grazie alla sua larghezza di banda e alla sua alta velocità di trasmissione che va da 1 Gb/s a 10Gb/s, molto più alta rispetto, ad esempio, a quella del CAN bus di 20Kbps.

La scelta di Automotive Ethernet è dovuta anche al crescente numero di sensori all'interno dei veicoli genera volumi significativi di dati, richiedendo potenti sistemi di elaborazione centralizzati. Per affrontare questo scenario, l'industria sta transitando verso un'architettura a zone, con veicoli dotati di massimo due computer ad alte prestazioni (il computer centrale e un altro che si occupa del backup

¹⁰ Gli Advanced Driver Assistance Systems (ADAS) sono tecnologie e sistemi progettati per migliorare la sicurezza e l'esperienza di guida attraverso l'automazione e l'assistenza nelle funzioni di guida. Questi sistemi utilizzano sensori, telecamere, radar, lidar e altri dispositivi per monitorare l'ambiente intorno al veicolo e fornire assistenza al conducente.

¹¹ Il termine "Zone ECU" fa riferimento a un Electronic Control Unit (ECU) configurato per gestire specifiche "zone" o aree funzionali all'interno di un veicolo o di un sistema complesso. Questo approccio suddivide le responsabilità di controllo tra diverse unità, ognuna dedicata a una funzione specifica come il motore, la trasmissione, la sicurezza o altre componenti. L'implementazione di Zone ECU consente una progettazione modulare e distribuita dei sistemi di controllo, facilitando la gestione e la manutenzione.

e all'occorrenza potrebbe sostituire il primo in caso di guasto) collegati in una rete ad anello insieme a tutti i sensori necessari.

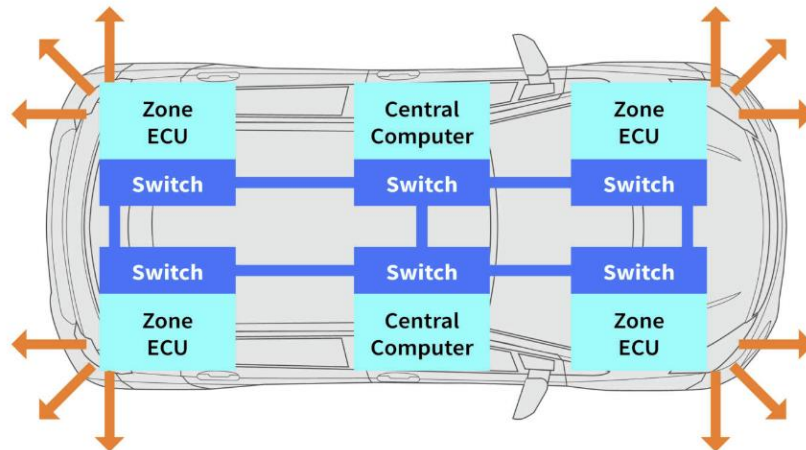


Figura 21: Connessione delle periferiche e del computer di bordo tramite Automotive Ethernet [62]

L'approccio a zone riduce il peso e la complessità del cablaggio, consentendo l'implementazione di funzioni come i fari abbaglianti adattivi senza richiedere una CU dedicata. La connettività è gestita da centraline di zona e HPC, riducendo la necessità di centraline locali complesse.

Tuttavia, le sfide emergono nella connettività dei System-on-Chip che, nonostante le elevate prestazioni di elaborazione, spesso mancano della connettività richiesta. Toshiba ha sviluppato un bridge di interfaccia di terza generazione, il TC9563XBG, che integra un'interfaccia PCIe Gen 3 con switch e una doppia interfaccia MAC TSN Ethernet in un package compatto. Questa soluzione offre connettività alle centraline di zona, alle piattaforme telematiche e ai sistemi di intrattenimento a bordo, affrontando le sfide di connettività attuali. Questo bridge si collega al SoC host tramite interfaccia PCIe, e fornisce la connettività Ethernet automotive fino a 10 Gb/sec . [62]

MEMORIA DEL COMPUTER DI BORDO DI UN' AUTO A GUIDA AUTONOMA

Analisi sulla miglior memoria per un'autovettura di livello 5

La grande mole di dati provenienti da tutti questi sensori deve essere gestita in maniera ottimale, in modo tale da poter garantire un flusso costante di dati ai sistemi di elaborazione; è essenziale dunque scegliere uno storage con capacità sufficiente e una velocità di lettura/scrittura adeguata a gestire questo flusso continuo.

La scelta di questo tipo di memoria deve essere un giusto mix tra larghezza di banda, capacità e potenza, ecco perché in questo capitolo andrò a confrontare le attuali migliori offerte, ossia le HBM2e, le RDIMM e le LPDDR5x.

Le HBM2e ossia le memorie a larghezza di banda elevata, sono un tipo di interfaccia di memoria utilizzata nelle DRAM ed offrono una larghezza di banda molto elevata pari a 410 GB/s, con una capacità massima di 24GB e una velocità massima di trasferimento pin a 3,2 Gbps con due canali bus da 128 bit per die. Tuttavia, l'adozione di HBM2e può presentare alcuni svantaggi. In primo luogo, i costi di produzione sono spesso più elevati rispetto a soluzioni di memoria tradizionali, influenzando il costo complessivo dei dispositivi. La complessità nella produzione, richiedendo processi avanzati, potrebbe impattare sulla disponibilità del prodotto sul mercato. Inoltre, la capacità massima potrebbe risultare leggermente inferiore rispetto ad alcune alternative, limitando la sua idoneità in contesti che richiedono grandi quantità di memoria.

Le Registered DIMM (RDIMM) costituiscono un tipo di modulo di memoria spesso impiegato in contesti server e sistemi ad elevate prestazioni di elaborazione dati. La caratteristica distintiva di queste memorie risiede nell'utilizzo di un registro o buffer aggiuntivo, atto a operare come intermediario tra il controller di memoria e i chip di memoria stessi. Il loro punto forte è sicuramente la capacità che può arrivare addirittura a 128 GB, ma a differenza delle HBM viste in precedenza le RDIMM hanno una più bassa larghezza di banda e anche una velocità di trasferimento inferiore.

In termini di affidabilità e stabilità il registro contribuisce a gestire i segnali di memoria, minimizzando il rischio di errori di lettura e scrittura e migliorando la qualità complessiva della connessione.

Le LPDDR5X invece sono un tipo di memoria DRAM progettata per dispositivi mobili come smartphone e tablet. Caratterizzate da velocità di trasferimento notevoli, offrono larghezze di banda elevate fino a 30 GB/s e sono ottimizzate per l'efficienza energetica. Queste memorie supportano una gamma variabile di capacità, solitamente comprese tra 4 GB e oltre 16 GB per modulo, in base alle specifiche del dispositivo in cui vengono utilizzate. Le LPDDR5X sono comunemente implementate in dispositivi mobili di fascia alta, contribuendo a garantire prestazioni avanzate e prolungare la durata della batteria grazie alla loro efficienza energetica.

Nel nostro caso, dunque, per lo storage di un computer di bordo per auto che deve implementare un IA in grado di guidare in maniera completamente autonoma, la scelta migliore potrebbe essere quella delle LPDDR, poiché garantiscono un consumo energetico nettamente inferiore oltre che dei costi decisamente più contenuti, ma anche per quanto riguarda la larghezza di banda, sebbene essa sia inferiore rispetto a quella delle HBM, , spesso forniscono prestazioni più che sufficienti per le esigenze specifiche delle applicazioni automotive. Questo è dimostrato dal fatto che anche le migliori offerte attualmente sul mercato come il Full Self-Driving (FSD) di Tesla e NVIDIA Drive Pegasus, che vedremo in seguito, fanno uso di memorie LPDDR.

Oltre a una memoria principale, tuttavia, il computer di bordo in questione dovrà avere a sua disposizione anche una memoria secondaria che permetta di mantenere i dati come mappe, dati di training, ecc..... in maniera permanente.

Per la memoria secondaria, una soluzione comune è l'utilizzo di unità di storage basate su tecnologie come SSD (Solid State Drive), in particolare l'utilizzo dello standard NVMe (Non-Volatile Memory Express) potrebbe essere una buona scelta poiché riduce fortemente i tempi di accesso ai dati fino a 3500 MB/s, con capacità di storage che arrivano anche a 2TB e sono anche conformi allo standard PCIe. [63]

UNITA' DI ELABORAZIONE DEL COMPUTER DI BORDO DI UN' AUTO A GUIDA AUTONOMA

Analisi sul miglior computer di bordo per un'IA di guida autonoma di
livello 5

Per effettuare l'inferenza di un'intelligenza artificiale in grado di pilotare autonomamente un veicolo, è essenziale disporre di un'unità di elaborazione potente che possa gestire carichi di lavoro intensi di deep learning in tempo reale. In generale, l'utilizzo di un System-on-a-Chip in un'automobile di livello 5, completamente autonomo, offre numerosi vantaggi. Questi includono un'efficace integrazione di componenti chiave, come CPU, GPU e acceleratori di intelligenza artificiale, massimizzando la potenza di calcolo in un'unica soluzione. Ciò migliora l'efficienza di integrazione e risparmia spazio, essenziale in un contesto automotive. Inoltre, la progettazione del SoC mira a garantire prestazioni elevate con un consumo energetico ottimizzato, critico per applicazioni in cui l'efficienza energetica è cruciale e la comunicazione ottimizzata tra i componenti sullo stesso chip contribuisce a ridurre la latenza.

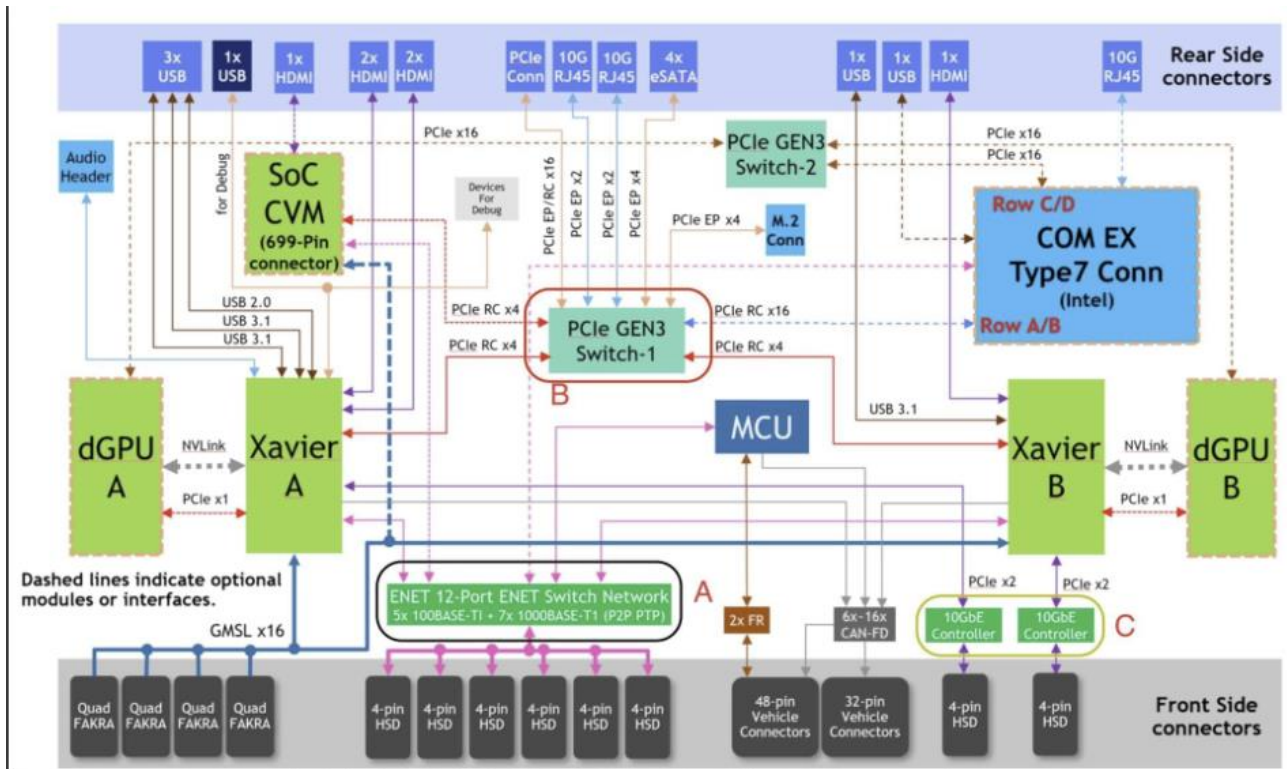


Figura 22: Schema a blocchi di NVIDIA DRIVE AGX Pegasus [66]

Un'opzione per un computer di bordo così potente potrebbe essere senza dubbio lo NVIDIA DRIVE AGX Pegasus, un System-on-a-Chip progettato per fornire elevate prestazioni per applicazioni di guida autonoma. Esso offre un'architettura combinata di due componenti principali: GPU basate sull'architettura NVIDIA Turing, tipiche delle GPU RTX da gaming (sono una valida scelta anche per quanto riguarda l'inferenza delle IA) e una coppia di SoC Xavier collegati grazie alla tecnologia NVIDIA NVLink per un totale di 320 TOPS [64].



Figura 23: Schema a blocchi SoC Xavier

Il SoC si compone di due unità Xavier, ciascuna delle quali include una CPU 8-core di fascia alta basata sull'architettura ARM v8 a 64 bit, divisa in 4 duplex ciascuno dei quali condivide 2MiB di cache L2 e una GPU con architettura Volta da 640 Tensor Core che dispone di otto multiprocessori stream Volta insieme ai loro 128 KiB standard di cache L1 e 512 KiB di L2 condivisa [66].

Oltre a CPU e GPU, l'NVIDIA Xavier è dotato anche di un Computer Vision Accelerator (CVA) e un'unità di elaborazione visiva (VP) ad alta gamma dinamica (HDR) da 8K e di 8 canali per le memorie LPDDR4X con una larghezza di banda massima di 127.1 GiB/s [67].

Il SoC integra anche una serie di sensori, tra cui telecamere, radar e altri dispositivi, che forniscono dati essenziali per la percezione dell'ambiente circostante da parte del veicolo autonomo. DRIVE PX Pegasus è progettato per gestire fino a 16 ingressi sensore dedicati ad alta velocità. Supporta anche Ethernet da 10 Gbit e ha una larghezza di banda di memoria di 1 Tbyte/s [65].

Un'altra opzione oltre quella di NVIDIA DRIVE AGX Pegasus è quella di Full Self-Driving Chip, un chip di guida autonoma progettato da Tesla per le proprie auto. Tesla afferma che il chip è destinato ai livelli autonomi 4 e 5. Il chip FSD incorpora 3 cluster quad-core Cortex-A72 per un totale di 12 CPU che funzionano a 2,2 GHz, una GPU Mali G71 MP12 che funziona a 1 GHz, 2 unità di elaborazione neurale che funzionano a 2 GHz e vari altri acceleratori hardware. L'FSD supporta una memoria LPDDR4-4266 fino a 128 bit [68].

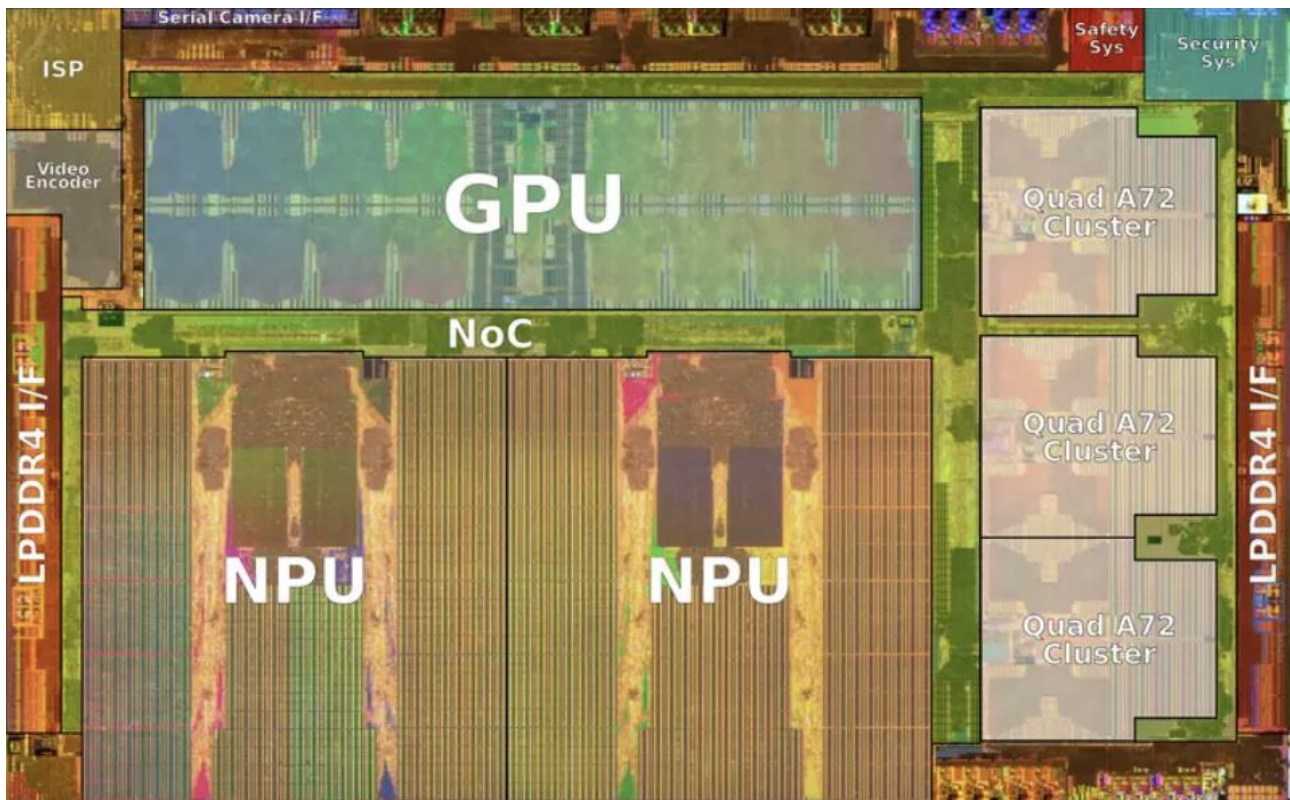


Figura 24: Full Self-Driving Chip di Tesla [68]

A livello generale, il chip rappresenta un sistema su chip completo in grado di avviare un sistema operativo standard. Il dispositivo è equipaggiato con dodici core ARM a 64 bit, organizzati in tre cluster di core quad-core Cortex-A72 che operano a 2,2 GHz per l'elaborazione generale. Inoltre, include una GPU relativamente leggera progettata principalmente per la post-elaborazione, funzionante a 1 GHz e con una capacità di fino a 600 GFLOPS, con supporto per operazioni in virgola mobile sia a precisione singola che doppia [68] [6].

Il chip FSD incorpora due unità di elaborazione neurale su misura, ciascuna con una memoria SRAM da 32 MiB progettata per archiviare temporaneamente i risultati di rete, riducendo la necessità di trasferire dati alla memoria principale. Operando a 2 GHz, ciascuna NPU ha una performance di picco di 36,86 trilioni di operazioni al secondo. Con due NPU su ciascun chip, il FSD può eseguire fino a 73,7 trilioni di operazioni al secondo di prestazioni di picco combinate [68].

La comparazione tra i chip FSD di Tesla e i chip Drive di NVIDIA evidenzia un importante trade-off tra prestazioni di picco ed efficienza energetica. Il chip FSD di Tesla offre 144 TOPS con un consumo di 72 watt, mentre il chip Pegasus di NVIDIA raggiunge i 320 TOPS ma richiede circa 500 watt.

Il vantaggio del chip Tesla, in termini di prestazioni per watt, suggerisce una maggiore efficienza energetica, il che può tradursi in una migliore gestione termica e potenzialmente in costi operativi inferiori. D'altra parte, i chip di nuova generazione di NVIDIA possono offrire prestazioni di picco superiori, ma a costo di un maggiore consumo energetico.

La scelta tra i due dipenderà dalle specifiche esigenze dell'applicazione e dalle restrizioni di potenza e termiche del veicolo.

CONCLUSIONI

In questa tesi, ho esplorato un argomento che ho immediatamente ritenuto affascinante e di estrema rilevanza. Il rapido sviluppo degli algoritmi e delle applicazioni di intelligenza artificiale sta plasmando in modo significativo il nostro presente e promette di rivoluzionare il nostro futuro. Questa evoluzione richiede risorse sempre più performanti, capaci di eseguire miliardi di operazioni al secondo, aprendo scenari inesplorati e stimolando la ricerca verso soluzioni architetture all'avanguardia.

Attualmente, l'attenzione degli ingegneri e dei programmatori è focalizzata sulla creazione di IA in grado di emulare il comportamento e il ragionamento umano. Tuttavia, le parole di Gary Marcus, professore alla New York University, ci spingono a considerare una prospettiva ancora più ambiziosa: *"La vera sfida dell'intelligenza artificiale non è replicare l'intelligenza umana, ma sviluppare forme di intelligenza che non esistono ancora."* Questa affermazione getta un'ombra di incertezza sul futuro delle architetture, poiché suggerisce che la vera innovazione non risiede solo nella replicazione di ciò che già conosciamo, ma nell'esplorazione di territori ancora inesplorati.

Il futuro delle architetture per l'IA rimane, dunque, un enigma affascinante. Questa tesi ha gettato le basi per comprendere le sfide e le opportunità attuali, ma lascia deliberatamente un punto interrogativo alla fine di questa conclusione. Come si evolveranno le architetture per soddisfare la sfida di sviluppare forme di intelligenza ancora sconosciute? Cosa ci riserverà il futuro dell'intelligenza artificiale? Solo il tempo e la continua ricerca porteranno risposte a questi interrogativi, aprendo le porte a un'avventura intellettuale che continua a sfidare e ispirare.

RINGRAZIAMENTI

Alla conclusione di questo elaborato, desidero dedicare un sentito ringraziamento a tutte le persone che hanno contribuito in modo significativo a questo meraviglioso percorso di approfondimento delle mie conoscenze durante gli anni universitari. Il mio sincero apprezzamento va alla mia Relatrice Diamantini e al mio correlatore Bordi, che con la loro disponibilità e gentilezza hanno guidato ogni passo di questo lavoro, dalla scelta dell'argomento fino alla sua conclusione.

Un grazie profondo va ai miei genitori, che con il loro costante supporto hanno reso possibile ogni conquista. A mia sorella Linda, un grazie speciale per aver reso il mio percorso più significativo con il suo affetto e incoraggiamento e a tutta la mia famiglia per il loro amore e supporto incondizionato.

Non posso dimenticare di ringraziare i miei amici di sempre per la loro costante presenza e il sostegno instancabile. Un ringraziamento speciale va anche ai miei compagni di università, con i quali ho condiviso gioie e sfide. Le esperienze condivise e il vostro supporto hanno reso questo viaggio accademico ancora più memorabile.

Questo percorso accademico è stato arricchito grazie a ognuno di voi.

Grazie di cuore.

FONTI BIBLIOGRAFICHE E SITOGRAFIA

Bibliografia

- [1] Gaurav Batra, Zach Jacobson, Siddarth Madhav, Andrea Queirolo, and Nick Santhanam. Artificial-intelligence hardware: New opportunities for semiconductor companies. McKinsey & Company. 2018.
- [2] Khan FH, Pasha MA, Masud S. Advancements in Microprocessor Architecture for Ubiquitous AI—An Overview on History, Evolution, and Upcoming Challenges in AI Implementation. *Micromachines*. 2021; 12(6):665.
- [3] Pudi Dhilleswararao, Srinivas Boppu, M. Sabarimalai Manikandan, Linga Reddy Cenkeramaddi. Efficient Hardware Architectures for Accelerating Deep Neural Networks: Survey. *IEEEAcces*. 2016; 4.
- [4] Maurizio Capra, Beatrice Bussolino, Alberto Marchisio, Guido Masera, Maurizio Martina, Muhammed Shafique. Hardware and Software Optimizations for Accelerating Deep Neural Networks: Survey of Current Trends, Challenges, and the Road Ahead. . *IEEEAcces*. 2020.
- [5] Zhixin Pan, Prabhat Mishra. Hardware Acceleration of Explainable Machine University of Florida, Department of Computer & Information Science & Engineering. 2022.

Sitografia

- [6] Toolify.ai Tesla FSD Chip: A Breakdown of Tesla vs Nvidia vs Intel Chips, <https://www.toolify.ai/gpts/tesla-fsd-chip-a-breakdown-of-tesla-vs-nvidia-vs-intel-chips-327588>. Consultato il 13/01/2024
- [7] Hewlett Packard Enterprise, Edge computing, <https://www.hpe.com/it/it/what-is/edge-computing.html#:~:text=L%27edge%20computing%20%C3%A8%20una,in%20un%20data%20center%20remoto>. Consultato il 3/12/2023.
- [8] TechCompany360, Intelligenza Artificiale: quali infrastrutture servono per abilitarla? <https://www.techcompany360.it/tech-lab/intelligenza-artificiale-quali-infrastrutture-servono-per-abilitarla/> . Consultato il 22/11/2023.
- [9] AI4Buisness, Sistemi IT: perché le infrastrutture per abilitare l'Intelligenza Artificiale devono evolvere, <https://www.ai4business.it/intelligenza-artificiale/intelligenzaartificiale-leinfrastruttureabilitanti/> . Consultato il 22/11/2023.
- [10] FastwebPlus, Cosa sono le TPU e quale il loro futuro, <https://www.fastweb.it/fastweb-plus/digital-magazine/cosa-sono-le-tpu-e-quale-il-loro-futuro/> . Consultato il 21/12/2023.
- [11] AI4Buisness, Intelligenza artificiale, quali sono le architetture hardware dedicate, <https://www.ai4business.it/intelligenza-artificiale/intelligenza-artificiale-quali-sono-le-architetture-hardware-dedicate/> . Consultato il 23/11/2023.
- [12] Medium, AI Optimised Hardware, <https://medium.com/appengine-ai/ai-optimised-hardware-da5e1d4c32b6> . Consultato il 29/11/2023.
- [13] HolisticSeo, AI Optimized Hardware: Definition, Importance, Examples and How It Works? <https://www.holisticseo.digital/ai/hardware/> . Consultato il 18/12/2023.
- [14] Thedecoder, What is AI hardware? Approaches, advantages and examples, <https://the-decoder.com/what-is-ai-hardware-approaches-advantages-and-examples/> . Consultato il 11/12/2023.
- [15] Deepgram, AI Hardware, <https://deepgram.com/ai-glossary/ai-hardware> . Consultato il 17/12/2023.
- [16] PugetSystems, Hardware Recommendations for MachineLearning/AI, <https://www.pugetsystems.com/solutions/scientific-computing->

- [workstations/machine-learning-ai/hardware-recommendations/](#) . Consultato il 4/12/2023.
- [17] C3.ai, What is Machine Learning? <https://c3.ai/introduction-what-is-machine-learning/> . Consultato il 02/12/2023.
 - [18] Tim Dettmers, Which GPU(s) to Get for Deep Learning: My Experience and Advice for Using GPUs in Deep Learning, <https://timdettmers.com/2023/01/30/which-gpu-for-deep-learning/> . Consultato il 5/12/2023.
 - [19] Medium, How do GPUs work? <https://medium.com/mlearning-ai/how-do-gpus-work-13bb243c17d> . Consultato il 5/12/2023.
 - [20] HDBlog, NVIDIA A100 AMPERE, <https://www.hdblog.it/hardware/articoli/n540322/nvidia-a100-80gb-acceleratore-gpu/#:~:text=Passando%20alla%20scheda%20tecnica%20della,tocca%20i%202%20TB%2Fs>. Consultato il 7/12/2023.
 - [21] Utmel, Neural Processing Unit NPU, Artificial Intelligence AI and Machine Learning ML explained, <https://www.utmel.com/blog/categories/integrated%20circuit/neural-processing-unit-npu-explained> . Consultato il 17/12/2023.
 - [22] BlackBlaze, AI 101: GPU vs TPU vs NPU, <https://www.backblaze.com/blog/ai-101-gpu-vs-tpu-vs-npu/> . Consultato il 18/12/2023.
 - [23] EDALAB, Machine Learning: introduzione agli algoritmi predittivi, <https://edalab.it/machine-learning-introduzione-agli-algoritmi-predittivi/> . Consultato il 2/12/2023.
 - [24] Dday, Nvidia ha il supercomputer per l'intelligenza artificiale più veloce del mondo: Perlmutter conta oltre 6.000 GPU, <https://www.dday.it/redazione/39662/nvidia-supercomputer-perlmutter-gpu> . Consultato il 14/12/2023.
 - [25] NVIDIA, GPU NVIDIA H100 Tensor Core, <https://www.nvidia.com/it-it/data-center/h100/> . Consultato il 14/12/2023.
 - [26] Moor, NVIDIA Targets Ampere Architecture To The Edge And 5G With EGX A100, <https://moorinsightsstrategy.com/nvidia-targets-ampere-architecture-to-the-edge-and-5g-with-egx-a100/> . Consultato il 28/12/2023.

- [27] TheNextPlatform, Diving Deep Into the NVIDIA Ampere GPU Architecture, <https://www.nextplatform.com/2020/05/28/diving-deep-into-the-nvidia-ampere-gpu-architecture/> . Consultato il 15/12/2023.
- [28] AnandTech, How GPUs work, <https://www.anandtech.com/show/7793/imaginations-powervr-rogue-architecture-exposed/2> . Consultato il 15/12/2023.
- [29] Techradar, What is an AI chip? Everything you need to know, <https://www.techradar.com/news/what-is-an-ai-chip-everything-you-need-to-know> . Consultato il 12/12/2023.
- [30] Intel, The Habana Gaudi2 Processor for Deep Learning, <https://www.intel.com/content/www/us/en/developer/articles/technical/habana-gaudi2-processor-for-deep-learning.html> . Consultato il 10/12/2023.
- [31] Nanoreview, Apple A17 Pro, <https://nanoreview.net/en/soc/apple-a17-pro> . Consultato il 15/12/2023.
- [32] AnandTech, NVIDIA Unveils “Titan RTX” Video Card: \$2500 Turing Tensor Terror Out Later This Month, <https://www.anandtech.com/show/13668/nvidia-unveils-rtx-titan-2500-top-turing> . Consultato il 18/12/2023.
- [33]Golden, NVIDIA H100 Tensor Core GPU, https://golden.com/wiki/NVIDIA_H100_Tensor_Core_GPU-5Z43AB5 . Consultato il 14/12/2023.
- [34] Sysgen, NVIDIA H100 TENSOR CORE GPU OVERVIEW, <https://www.sysgen.de/en/blog/hardware/introducing-nvidia-hgx-h100-an-accelerated-server-platform-for-ai-and-high-performance-computing> . Consultato il 14/12/2023.
- [35] DeepGram, Tensor Processing Unit(TPU), <https://deepgram.com/ai-glossary/tensor-processing-unit-tpu> . Consultato il 03/01/2024.
- [36] Vyrian, What is a Tensor Processing Unit and How does it Work? <https://www.vyrian.com/what-is-a-tensor-processing-unit-tpu-and-how-does-it-work/> . Consultato il 02/01/2024.
- [37] Medium, AI Chips: Google TPU, <https://jonathan-hui.medium.com/ai-chips-tpu-3fa0b2451a2d> . Consultato il 22/12/2023.
- [38] Medium, What’s inside a TPU? <https://medium.com/@antonpaquin/whats-inside-a-tpu-c013eb51973e> . Consultato il 22/12/2023.

- [39] Google, Architettura di sistema, <https://cloud.google.com/tpu/docs/system-architecture-tpu-vm?hl=it> . Consultato il 22/12/2023.
- [40] DataCenterDynamics, AWS makes Trainium EC2 instances generally available, <https://www.datacenterdynamics.com/en/news/aws-makes-trainium-ec2-instances-generally-available/> . Consultato il 04/01/2024.
- [41] AWS, NeuronCore-v2 Architecture, <https://awsdocs-neuron.readthedocs-hosted.com/en/latest/general/arch/neuron-hardware/neuron-core-v2.html#neuroncores-v2-arch> . Consultato il 04/01/2024.
- [42] AWS, Trainium Architecture, <https://awsdocs-neuron.readthedocs-hosted.com/en/latest/general/arch/neuron-hardware/trainium.html> . Consultato il 03/01/2024.
- [43] AWS, AWS Trn1/Trn1n Architecture, <https://awsdocs-neuron.readthedocs-hosted.com/en/latest/general/arch/neuron-hardware/trn1-arch.html> . Consultato il 04/01/2024.
- [44] AWS, AWS Inf2 Architecture, <https://awsdocs-neuron.readthedocs-hosted.com/en/latest/general/arch/neuron-hardware/inf2-arch.html> . Consultato il 04/01/2024.
- [45] AWS, Inferentia2 Architecture, <https://awsdocs-neuron.readthedocs-hosted.com/en/latest/general/arch/neuron-hardware/inferentia2.html> . Consultato il 04/01/2024.
- [46] AWS, Istanze DL1 di Amazon EC2, <https://aws.amazon.com/it/ec2/instance-types/dl1/> . Consultato il 19/12/2023.
- [47] AWS, Istanze P4 di Amazon EC2, <https://aws.amazon.com/it/ec2/instance-types/p4/> . Consultato il 19/12/2023.
- [48] AWS, Amazon EC2 P4d instances deep dive, <https://aws.amazon.com/it/blogs/compute/amazon-ec2-p4d-instances-deep-dive/> . Consultato il 19/12/2023.
- [49] NetworkDigital365, Deep learning: cos'è, come funziona e applicazioni - Agenda Digitale, <https://www.agendadigitale.eu/cultura-digitale/deep-learning-cos-come-funziona-e-applicazioni/> . Consultato il 2/12/2023.
- [50] Smart, Algebra lineare e sue applicazioni al Machine Learning, Deep Learning e Big Data Analysis, <https://smartstrategy.eu/intelligenza-artificiale/algebra-lineare-e-sue-applicazioni-al-machine-learning-deep-learning-e-big-data-analysis/> . Consultato il 15/12/2023.

- [51] Synopsys, The DNA of an Artificial Intelligence SoC, https://www.synopsys.com/designware-ip/technical-bulletin/the-dna-of-an-ai-soc-dwtb_q318.html . Consultato il 12/12/2023.
- [52] DzTechs, Cos'è il Neural Engine di Apple e come funziona? <https://www.dztechs.com/it/what-is-a-neural-engine-how-does-it-work/> . Consultato il 15/12/2023.
- [53] Habana, Efficient scale-out for the large-scale generative AI era. <https://habana.ai/products/networking/> . Consultato il 10/12/2023.
- [54] ElectronicsHub, CUDA cores to accelerate your workflows, <https://www.electronicshub.org/cuda-cores/> . Consultato il 15/12/2023.
- [55] PC HI-TECH, NVIDIA TITAN RTX, il dinosauro delle GPU, <https://www.pcprofessionale.it/news/nvidia-titan-rtx-dinosauro-gpu/> . Consultato il 20/12/2023.
- [56] CardinalPeak, At The Edge Vs. In The Cloud: Artificial Intelligence And Machine Learning, <https://www.cardinalpeak.com/blog/at-the-edge-vs-in-the-cloud-artificial-intelligence-and-machine-learning> . Consultato il 3/12/2023.
- [57] NetworkDigital360, AI e cloud computing: vantaggi e svantaggi della combinazine, <https://www.zerounoweb.it/cloud-computing/ai-e-cloud-computing-vantaggi-e-svantaggi-della-combinazione/> . Consultato il 3/12/2023.
- [58] NVIDIA, NVIDIA Introduces DRIVE AGX Orin, <https://nvidianews.nvidia.com/news/nvidia-introduces-drive-agx-orin-advanced-software-defined-platform-for-autonomous-machines> . Consultato il 7/01/2024.
- [59] Arrow, Addestramento dei veicoli autonomi e del motore dati di Tesla, <https://www.arrow.com/it-it/research-and-events/articles/autonomous-vehicle-training-and-teslas-data-engine-explained> . Consultato il 7/01/2024.
- [60] Forbes, New MLPerf Benchmarks Show Why NVIDIA Remorked Its Product Roadmap, <https://www.forbes.com/sites/karlfreund/2023/11/08/new-mlperf-benchmarks-show-why-nvidia-has-reworked-its-product-roadmap/?sh=33b864bb655c> . Consultato il 10/01/2024.
- [61] Fare Elettronica, Automobili a Guida Autonoma: Come funzionano e i diversi Livelli, <https://fareelettronica.it/automobili-guida-autonoma/#:~:text=Le%20automobili%20a%20guida%20autonoma%20sono%20dotate%20di%20una%20serie,raccolgono%20dati%20dall'ambiente%20circostante> . Consultato il 7/01/2024.

- [62] elettronicanews, Tecnologia Ethernet nelle architetture a zone, <https://www.elettronicanews.it/tecnologia-ethernet-nelle-architetture-a-zone/> . Consultato l'8/01/2024.
- [63] tom'sHARDWARE, what are HBM, HBM2 and HBM2E? <https://www.tomshardware.com/reviews/glossary-hbm-hbm2-high-bandwidth-memory-definition,5889.html> . Consultato il 10/01/2024.
- [64] NVIDIA, Soluzioni di guida autonoma e sicura, <https://www.nvidia.com/it-it/deep-learning-ai/solutions/inference-platform/automotive/> . Consultato il 11/01/2024.
- [65] EletronicDesign, NVIDIA?s Volta Architecture Gives DRIVE PX Pegasus Its Smarts, <https://www.electronicdesign.com/markets/automotive/article/21805735/nvidias-volta-architecture-gives-drive-px-pegasus-its-smarts> . Consultato il 11/01/2024.
- [66] NVIDIA Activating 10 GbE Controllers in Pegasus Board, <https://forums.developer.nvidia.com/t/activating-10gbe-controllers-in-pegasus-board/190058> . Consultato il 12/01/2024.
- [67] WikiChip, Tegra Xavier – Nvidia, <https://en.wikichip.org/wiki/nvidia/tegra/xavier> . Consultato il 12/01/2024.
- [68] WikiChip, FSD Chip – Tesla, [https://en.wikichip.org/wiki/tesla_\(car_company\)/fsd_chip](https://en.wikichip.org/wiki/tesla_(car_company)/fsd_chip) . Consultato il 13/01/2024.
- [69] GizChina, Apple A17 Pro ufficiale: è il primo SoC a 3 nm della storia, <https://gizchina.it/2023/09/apple-a17-pro-ufficiale-novita/> . Consultato il 20/12/2024.
- [70] Intel, Il processore Habana Gaudi2 per il deep learning, <https://www.intel.cn/content/www/cn/zh/developer/articles/technical/habana-gaudi2-processor-for-deep-learning.html> . Consultato il 23/12/2023.
- [71] STH, Intel Gaudi to Gaudi2 Specs, <https://www.servethehome.com/intel-habana-gaudi2-launched-ai-training-chip-supermicro-ddn-oam/intel-gaudi-to-gaudi2-specs/> . Consultato il 23/12/2023.
- [72] Habana, Efficient scale-out for the large-scale generative AI era, <https://habana.ai/products/networking/> . Consultato il 20/12/2023.
- [73] NVIDIA, The Ultimate PC GPU NVIDIA TITAN RTX, <https://www.nvidia.com/content/dam/en-zz/Solutions/titan/documents/titan-rtx-for-creators-us-nvidia-1011126-r6-web.pdf> . Consultato il 18/12/2023.

- [74] naddod, Optical Transceivers and Artificial Intelligence: Jointly Creating an Intelligent Era, <https://www.naddod.com/blog/optical-transceivers-and-artificial-intelligence-jointly-creating-an-intelligent-era> . Consultato il 16/12/2023.
- [75] Medium, DGX H100 Super Computer, <https://medium.com/@neuralinternet/dgx-h100-super-computer-754b86e852ae> . Consultato il 7/01/2024.
- [76] AnandTech, NVIDIA Gives Xavier Status Update & Announces TensorRT 3 at GTC China 2017 Keynote, <https://www.anandtech.com/show/11872/nvidia-xavier-status-update-and-tensorrt-3-announcement-at-gtc-china-2017-keynote> . Consultato il 12/01/2024.
- [77] AI4Buisness, Intelligenza Artificiale, cos'è, come funziona, le applicazioni. <https://www.ai4business.it/intelligenza-artificiale/intelligenza-artificiale-cose/> . Consultato il 23/11/2023.