

Vettorizzazione di operazioni aritmetiche a 4 bit utilizzando registri a 64 bit sull'architettura x86-64

Tesi di Laurea in Ingegneria Informatica

Candidato

Lorenzo Grassi

Relatore

Prof. Marco Cococcioni



UNIVERSITÀ DI PISA

Introduzione e Problema

- I *Large Language Model* utilizzano miliardi di parametri. Per risparmiare spazio, questi vengono a volte memorizzati su pochissimi bit.
- I processori non hanno istruzioni che lavorano direttamente su tipi più piccoli di 1 byte.
- Utilizzando tipi dato a 4 bit, possiamo inserire 16 valori in ogni registro *general-purpose*. In particolare usiamo interi *unsigned*.
- Il nostro obiettivo è eseguire operazioni in parallelo tra i *nibble* corrispondenti di due registri, in pratica facendo in software quello che fanno le istruzioni SIMD.

0001	0110	0010	1100	1101	0100	0101	1010	0101	0110	1111	0110	0100	0100	0011	1100
------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------

Registro a 64 bit interpretato come 16 valori a 4 bit.

- Utilizzando operazioni aritmetiche e bit a bit, abbiamo implementato varie procedure che calcolano somma, differenza, e prodotto dei *nibble* corrispondenti in modo vettorizzato. I 16 risultati sono anch'essi su *nibble* di un registro a 64 bit.
- Abbiamo scritto anche le varianti con saturazione, che in caso di overflow limitano il risultato al massimo (o minimo) possibile.

			11	1111111	1111			11	111	1							
1001	0000	0110	0001	0001	0011	1101	1010	0001	0011	1110	0101	1011	1110	0100	0000		+
0010	0100	0111	0110	1111	1000	0110	0000	1011	0010	1000	0110	0000	0000	1011	1010		=
1011	0100	1101	1000	0000	1100	0011	1010	1100	0110	0110	1011	1011	1110	1111	1010		
0000	0000	1100	1111	1110	0111	1000	0000	0110	0111	0000	1000	0000	0000	0000	0000		

			111		11				1								
1011	0100	1101	1000	0000	1100	0011	1010	1100	0110	0110	1011	1011	1110	1111	1010		-
0000	0000	0000	0001	0000	0001	0000	0000	0000	0001	0000	0000	0000	0000	0000	0000		=
1011	0100	1101	0111	0000	1011	0011	1010	1100	0101	0110	1011	1011	1110	1111	1010		

Esempio di esecuzione dell'algoritmo a doppio XOR per la somma. Questa esecuzione ha richiesto due passaggi.

- Oltre alle semplici operazioni aritmetiche, abbiamo anche implementato operazioni più complesse dell'algebra lineare, in particolare:
 - Interpretando i due registri di input come vettori a 16 componenti, abbiamo scritto un algoritmo per calcolare il prodotto scalare.
 - Abbiamo implementato la moltiplicazione e accumulo su corsia (*multiply-accumulate lane*), che prende tre registri (a, b, c) e uno scalare (lane) come input (ciascuno come vettore di 16 valori a 4 bit) e, per ogni componente i , calcola $a[i] + b[i] * c[lane]$.
 - Infine, sfruttando quest'ultima, abbiamo scritto una procedura che calcola il prodotto tra matrici con elementi a 4 bit in modo vettorizzato.
 - Per queste ultime due operazioni abbiamo progettato anche una versione con saturazione.

- Abbiamo confrontato queste operazioni vettorizzate con alternative non vettorizzate, e con funzioni basate su *lookup table*.
- La *lookup table* è risultata quasi sempre più lenta, tranne che per la moltiplicazione con saturazione.
- Rispetto all'approccio non vettorizzato, l'approccio vettorizzato è risultato tra le 3 e le 5 volte più veloce per quasi tutte le operazioni, e quasi 7 volte più veloce per la moltiplicazione tra matrici.
- Compilando con -O3, GCC ha vettorizzato il codice inizialmente scalare di alcune operazioni, utilizzando veri e propri registri SIMD. Queste ottimizzazioni competono e in molti casi superano le nostre operazioni vettorizzate su registri *general-purpose*.