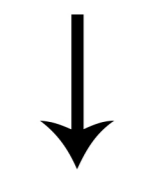
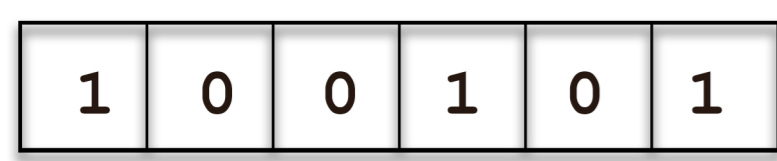


Revisiting POPCOUNT Operations in CPUs/GPUs

Chenfan Sun - chenfs@cs.washington.edu

Motivation

POPCOUNT:



3

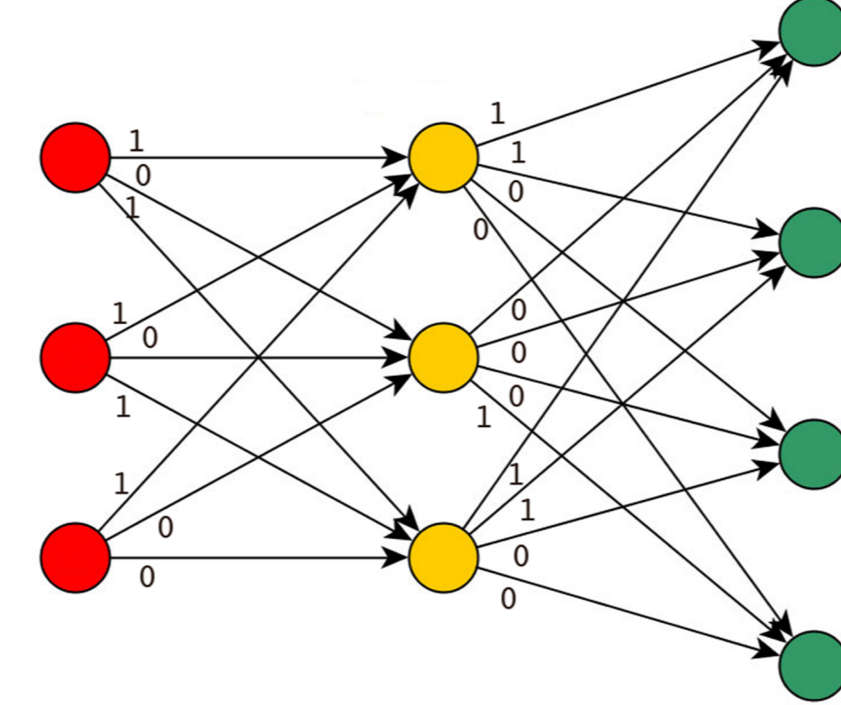
Hamming Distance[1]

$$D_H = \sum_{i=1}^k |x_i - y_i|$$

$x = y \Rightarrow D = 0$

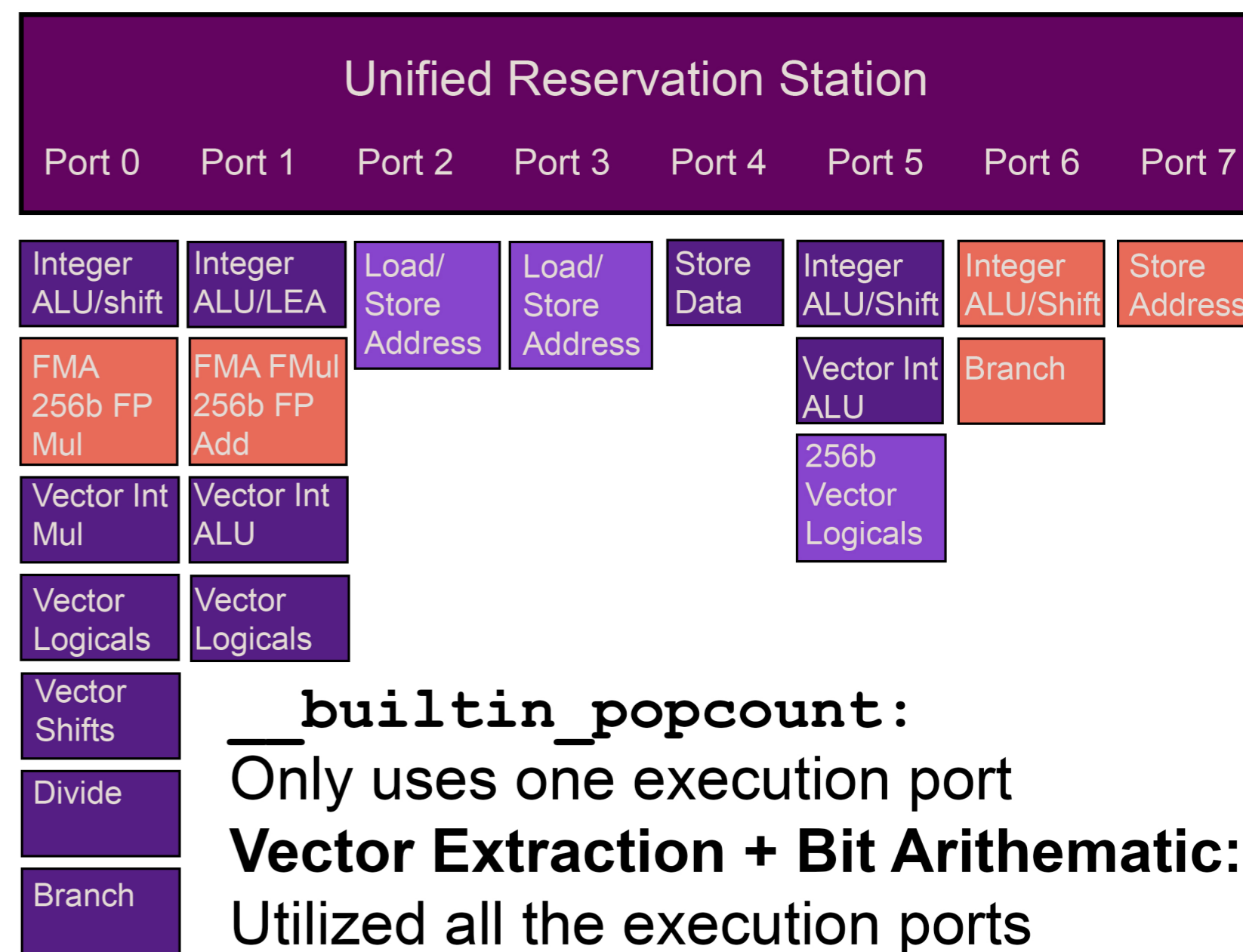
$x \neq y \Rightarrow D = 1$

Binary Neural Networks [2]



Background

Intel Haswell Execution Ports



`__builtin_popcount:`

Only uses one execution port

Vector Extraction + Bit Arithmetic:
Utilized all the execution ports

Approach

Data Set Size: 2MB (Fits in Cache), 1GB

CPUs

Vector Extension

{ SSSE3
AVX-2 }

Algorithm

{ Harley-Seal
Swar
Lookup }

Testbed

Intel Haswell i7-4790K
GeForce GTX Titan X (Maxwell)
DDR3 1333MHz Dual-Channel

GPUs

Data Width

{ 32 bits
64 bits
128 bits }

Algorithm

{ CUDA popcnt
Swar
Harley-Seal }

Reduction

{ CPU Reduction
GPU Reduction
GPU Shuffle }

```
a = 0b1101 # initial number
a0 = (a >> 0) & 0b0101 # every other bits
a1 = (a >> 1) & 0b0101 # remaining bits
b = a0 + a1 # 1s in 2-bit slice
b0 = (b >> 0) & 0b0011 # = 0b0001
b1 = (b >> 2) & 0b0011 # = 0b0010
ans = b0 + b1 # 1s in 4-bit slice,
ans: 3 # final answer in 2-bit word
```

Conclusion

Better Performance for *Compute Bound*:

More Powerful Processor



Haswell
i7-4790K



Haswell
i7-6700K

Wider Vector Extension

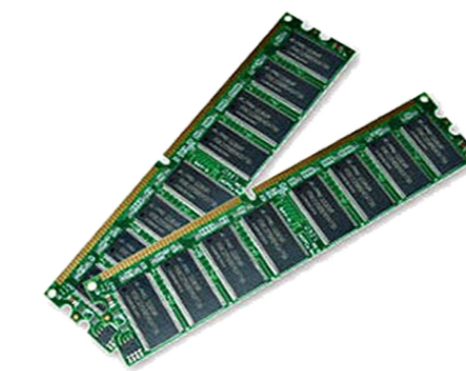
{ SSSE3
AVX-2 } \Rightarrow { AVX-512 }

Better Algorithm

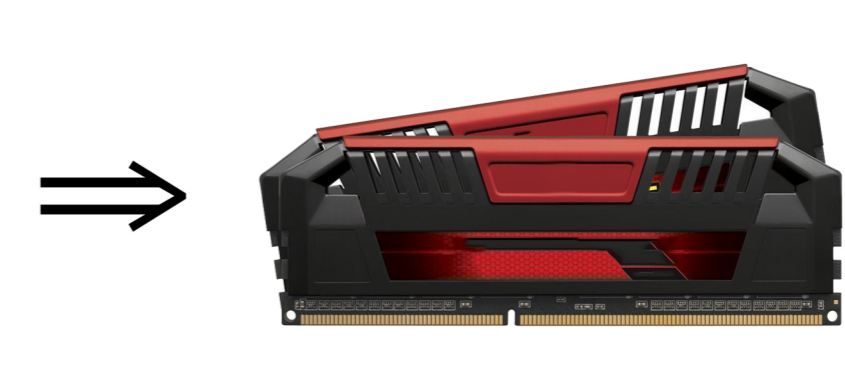
{ Lookup
Table } \Rightarrow { Harley-Seal }

Better Performance for *Memory Bound*:

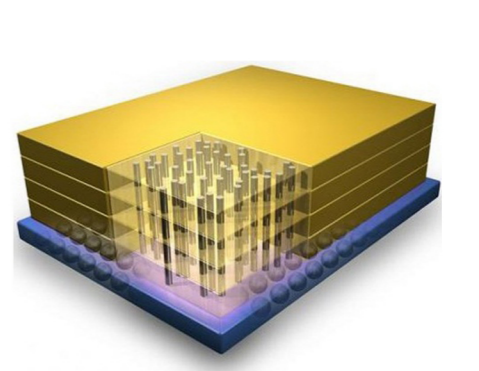
Wider memory channel / Faster clock rate / Architecture



Dual-Channel
DDR3 1333MHz



Quad-Channel
DDR4 2133MHz



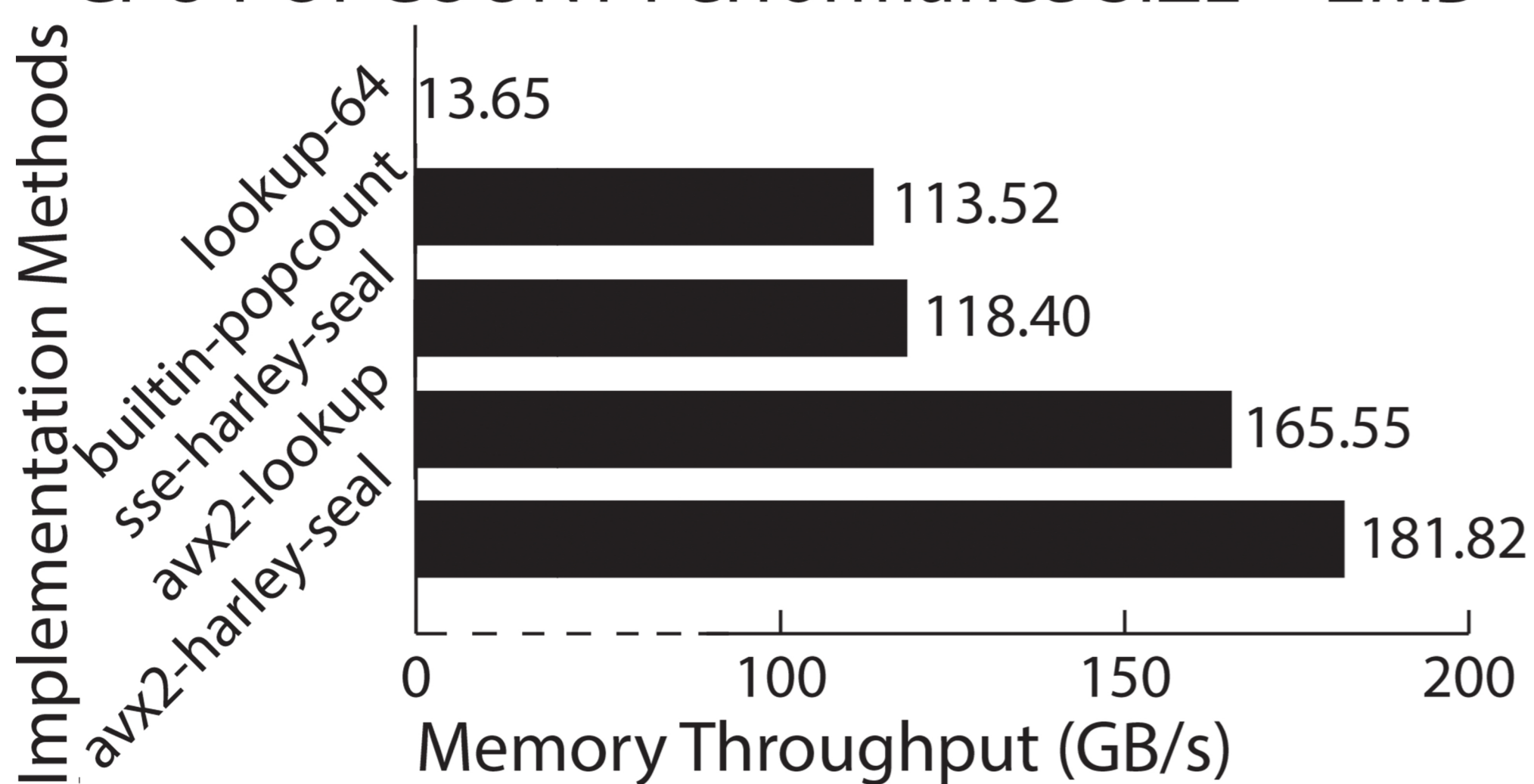
Die-Stack
Memory

Algorithmic choices do not matter for memory bound.

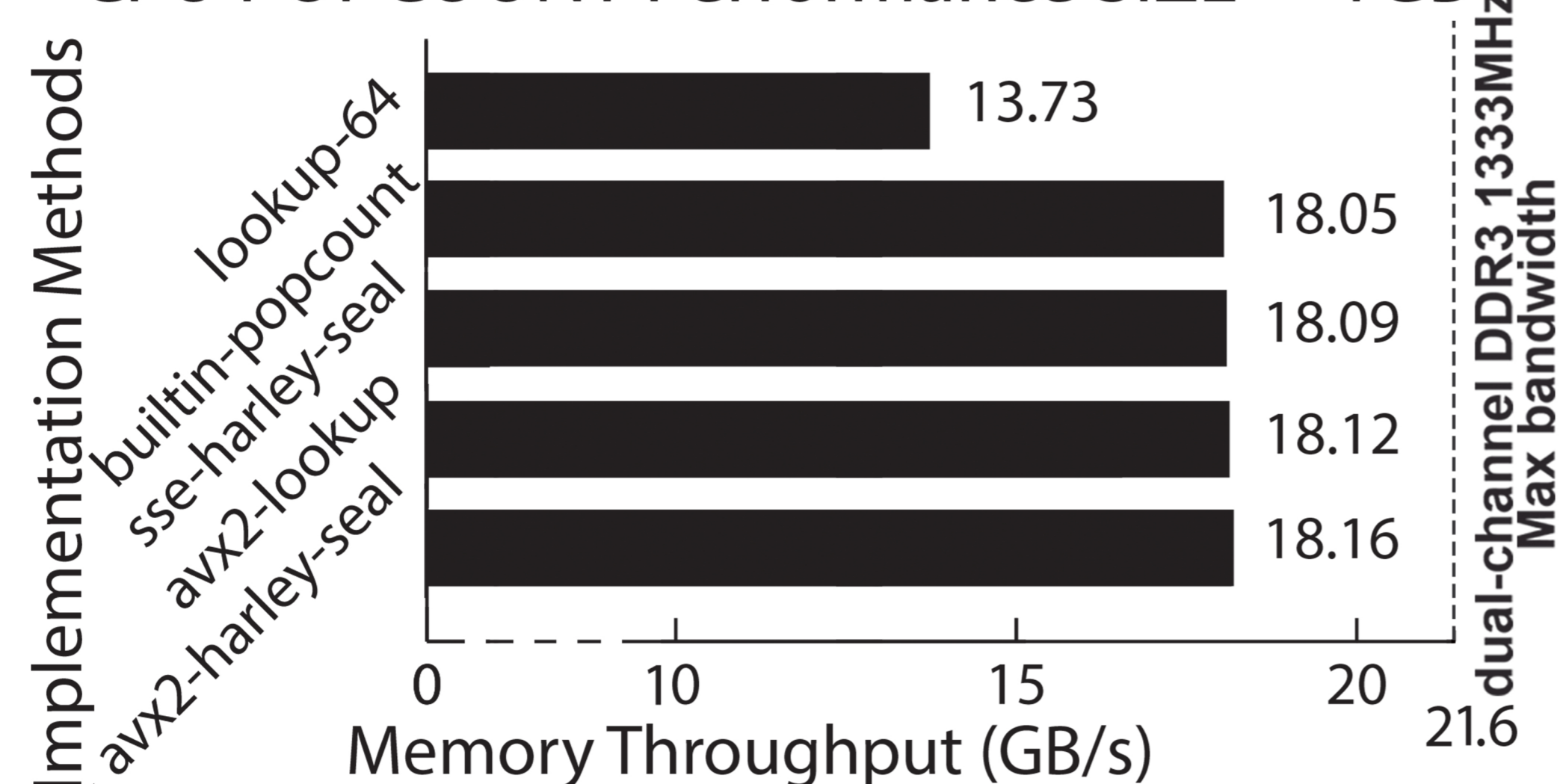
Results

CPUs

CPU POPCOUNT Performance SIZE = 2MB

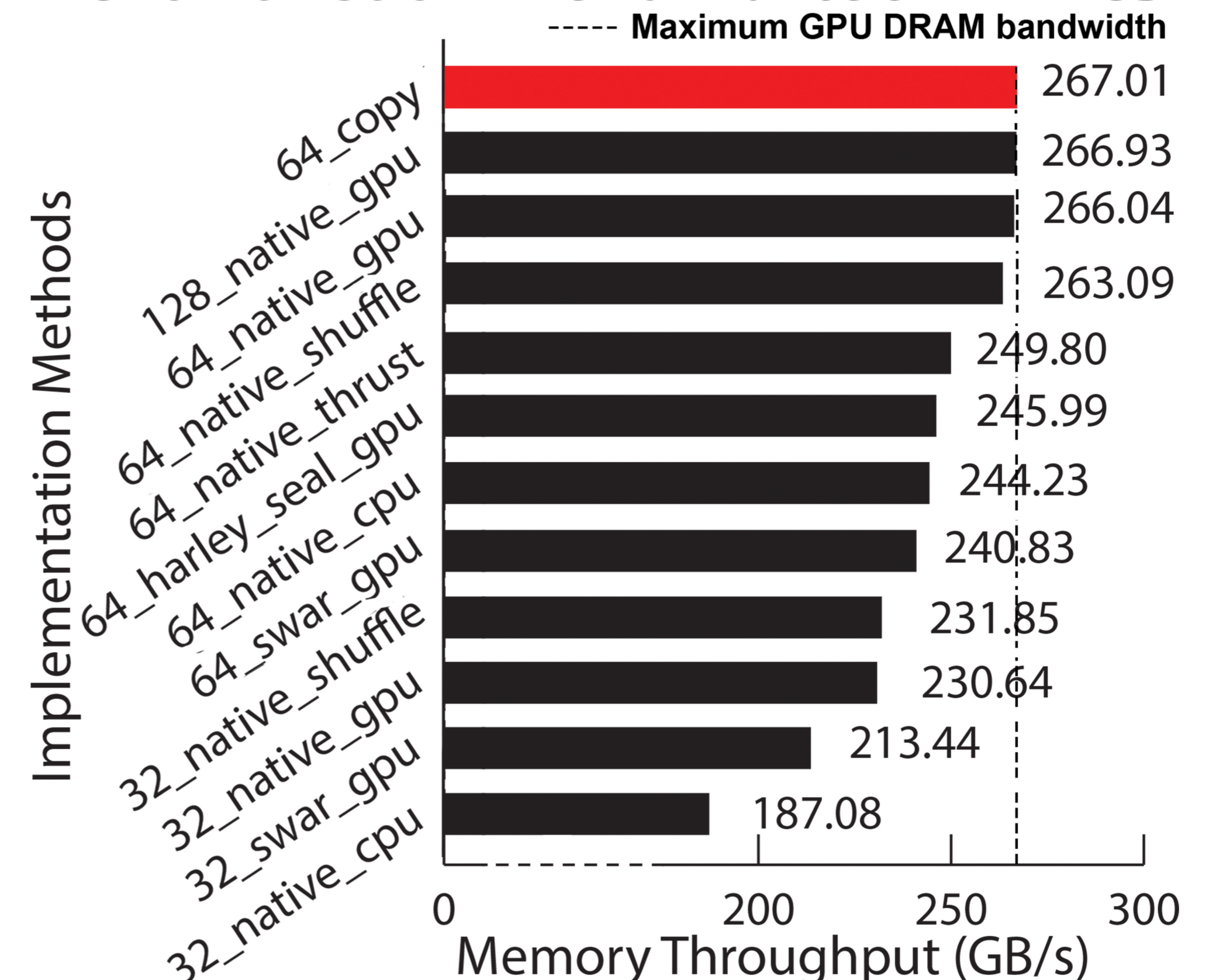


CPU POPCOUNT Performance SIZE = 1GB



GPUs

GPU POPCOUNT Performance SIZE = 1GB



[1] V.T. Lee et al., "NCAM: Near-Data Processing for Nearest Neighbor Search," in *arXiv 2016*.

[2] M. Rastegari et al., "XNOR-Net: ImageNet Classification Using Binary Convolutional Neural Networks," in *ECCV 2016*, Amsterdam, Netherlands.