Disclaimer: REBE The contents of this talk do not necessarily reflect the quality, rigor, and integrity of my lab's actual technical work.





hpcaaraae.ora/salishan

Neural networks,





















VIVA LA EVOLUCIÓN





hpcgarage.org/salishan

Neural networks, human brains, and "The Memoristas"

Richard (Rich) Vuduc









ImageNet Dataset

Question: How much energy does learning require?

Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., ... & Fei-Fei, L. (2015). Imagenet large scale visual recognition challenge. arXiv preprint arXiv:1409.0575. [web]

IM A GENET





ImageNet Dataset



Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., ... & Fei-Fei, L. (2015). Imagenet large scale visual recognition challenge. arXiv preprint arXiv:1409.0575. [web]

IM **GENET**





(4 GPUs) x (250 Watts / GPU) x (1 week) ~ 0.6 billion Joules

hpcgarage.org/salishan





(4 GPUs) x (250 Watts / GPU) x (1 week) ~ 0.6 billion Joules

(1 brain) x (1 year)

x (20 Watts / brain) ~ 0.6 billion Joules 71223832

gettyimages* Roderick Chen





(4 GPUs) x (250 Watts / GPU) x (1 week) ~ 0.6 billion Joules

22

(1 brain) x (1 year)

x (20 Watts / brain) ~ 0.6 billion Joules 71223832 gettyimages* Roderick Chen



The Iron Law of parallel and distributed communication costs







Two costs: *T*_{network} + *T*_{memory}





The Iron Law of parallel and distributed communication costs



hpcgarage.org/salishan





The ron Law of parallel and distributed communication costs



hpcgarage.org/salishan

B = "balance" - Kung (1986); Hillis (1987); Callahan & Kennedy (1988); Blelloch, Maggs, Miller (1994); McCalpin (1995); Williams et al. (2009)See also: Czechowski et al. (ICS'2012); Young & Vuduc (PMES'2016)

Two costs: *T*_{network} + *T*_{memory}

A fixed comm-bound algorithm + fixed problem size + fixed machine peak:

$$\frac{\beta_{\rm mem}^{\kappa}}{\beta_{\rm link}} \bigg) \,, \qquad \kappa = (d-1)/d \qquad \text{(d-dimensional torus)}$$

Bandwidth (GB/s)

Bandwidth (GB/s)

Bandwidth (GB/s)

$\sim 10^{-3} \sim 10^{10^{-3}}$

http://www.jetpress.org/volume1/moravec.pdf

Communication capability dominates that of computation – huh?

This ratio almost appears nonsensical as we think of computation today. Is it?

Interpretations: New computational models are needed, perhaps that reduce, sample, and approximate judiciously, or where the mere act of moving bits is the computation.

 $10^2 -$

10⁰ ->

Jerformance

10⁻⁴ -

single GFLOP/s)

Communication capability dominates that of computation — huh?

This ratio almost *appears* nonsensical as we think of computation **today**. Is it?

Interpretations: New computational models are needed, perhaps that **reduce**, **sample**, and **approximate** judiciously, or where the mere act of moving bits **is** the computation.

(Single GFLOP/s) 10⁻² -10⁻⁴ -

 $10^2 -$

Rich Vuduc

High-performance computing, scalable parallel algorithms & software

College of Computing **Computational Science and Engineering**

hpcgarage.org/salishan