# How to Write Fast Numerical Code

Fall 2016
*Lecture:* Balance Principles, Part II

**Instructor:** Torsten Hoefler & Markus Püschel

**TA:** Salvatore Di Girolamo

**ETH**
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

# References

- **These slides and the work is from Kenneth Czechowksi, Rich Vuduc et al., Georgia Tech**

- Kenneth Czechowski, Casey Battaglino, Chris McClanahan, Aparna Chandramowlishwaran, and Richard Vuduc. **Balance principles for algorithm-architecture co-design.** In *Proc. USENIX Wkshp. Hot Topics in Parallelism (HotPar)*, May 2011.

- Kenneth Czechowski, Chris McClanahan, Casey Battaglino, Kartik Iyer, P.-K. Yeung, Richard Vuduc. **On the communication complexity of 3D FFTs and its implications for exascale.** In *Proceedings of the ACM International Conference on Supercomputing (ICS)*, 2012.

# Balance Principles II

*Czechowksi et al. 2011*

$$T_{\text{mem}} \leq T_{\text{comp}}$$

$$\frac{p\pi}{\beta}\left(1 + \frac{\alpha\beta/\lambda}{Q/D}\right) \leq \frac{W}{Q\lambda}\left(1 + \frac{p}{W/D}\right)$$

# Application: Analyze Effect of HW Trends

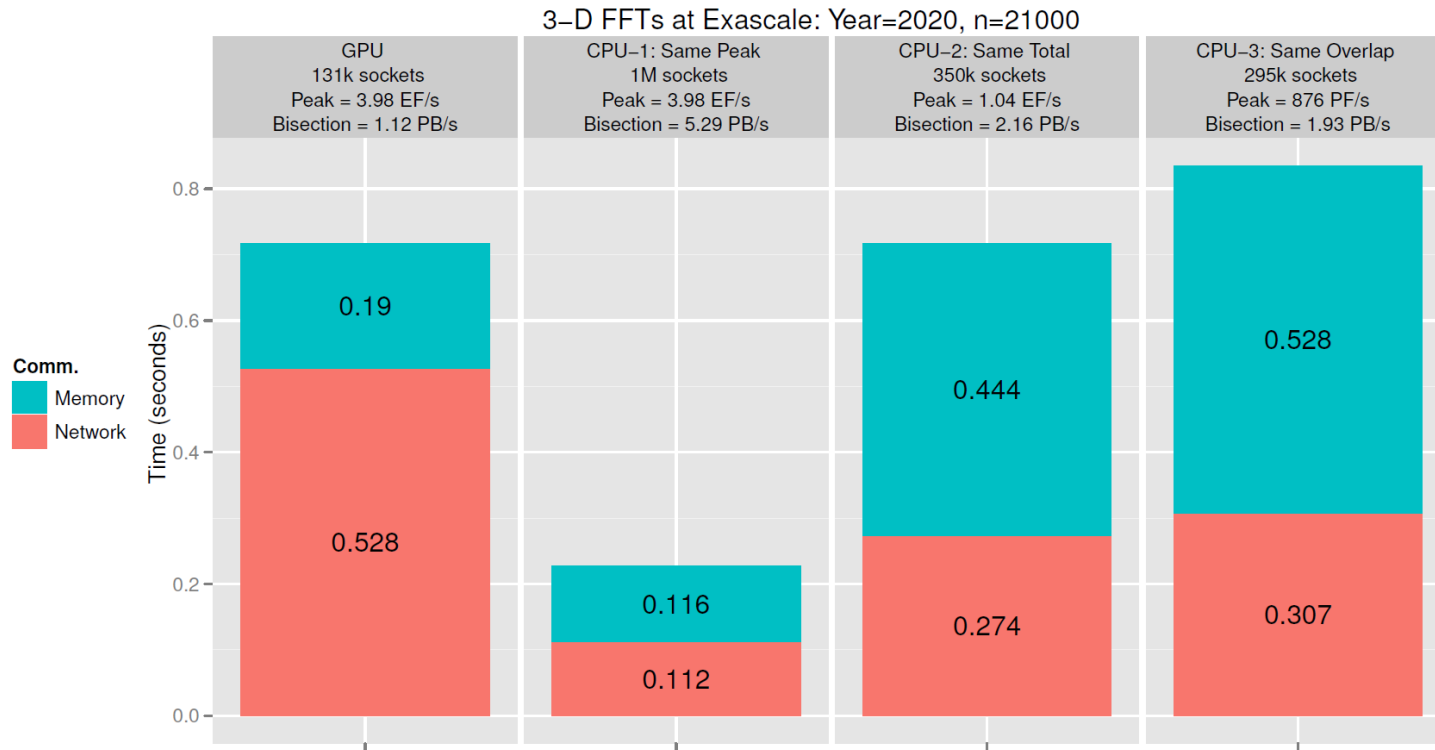*Czechowksi et al. 2012*

## 10 year extrapolation (2010 – 2020)

| Parameter | | 2010 values | Doubling time (in years) | 10-year increase factor | value |
|---|---|---|---|---|---|
| Peak: | $C_{\mathrm{CPU}}$ | 50.4 GF/s | 1.7 | 59.0× | 3.0 TF/s |
| | $C_{\mathrm{GPU}}$ | 515 GF/s | | | 30 TF/s |
| Cores:[a] | $\rho_{\mathrm{CPU}}$ | 6 | 1.87 | 40.7× | 134 |
| | $\rho_{\mathrm{GPU}}$ | 448 | | | 18k |
| Memory bandwidth: | $\beta_{\mathrm{CPU}}$ | 21.3 GB/s | 3.0 | 9.7× | 206 GB/s |
| | $\beta_{\mathrm{GPU}}$ | 144 GB/s | | | 1.4 TB/s |
| Fast memory | $Z_{\mathrm{CPU}}$ | 6 MB | 2.0 | 32.0× | 192 MB |
| | $Z_{\mathrm{GPU}}$ | 2.7 MB[b] | | | 86.4 MB |
| Line size: | $L_{\mathrm{CPU}}$ | 64 B | 10.2 | 2.0× | 128 B |
| | $L_{\mathrm{GPU}}$ | 128 B | | | 256 B |
| Link bandwidth: | $\beta_{\mathrm{link}}$ | 10 GB/s | 2.25 | 21.8× | 218 GB/s |
| Machine peak: | $R_{\mathrm{peak}}$ | 4 PF/s | 1.0 | 1000× | 4 EF/s |
| System memory: | $E$ | 635 TB | 1.3 | 208× | 132 PB |
| Nodes ($\frac{R_{\mathrm{peak}}}{C}$): | $P_{\mathrm{CPU}}$ | 79,400 | 2.4 | 17.4× | 1.3M |
| | $P_{\mathrm{GPU}}$ | 7,770 | | | 135,000 |

# Application: Analyze Effect of HW Trends

*Czechowksi et al. 2012*

**3D-FFT in 2020:**

**Faster on CPU or GPU?**



3−D FFTs at Exascale: Year=2020, n=21000
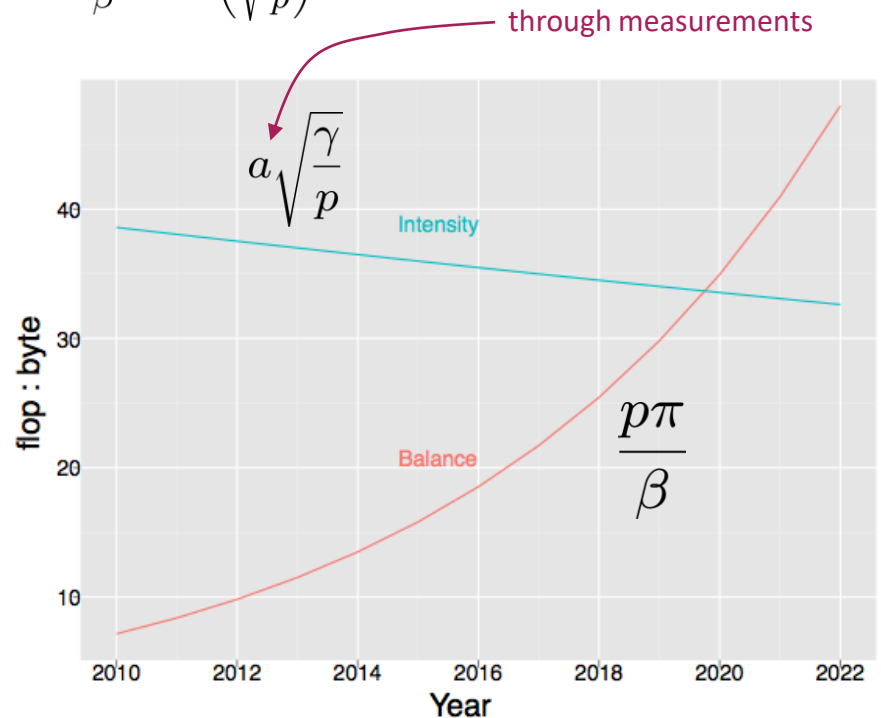
# Application: Analyze Effect of HW Trends

*Czechowksi et al. 2012*

## 10 year extrapolation (2010 – 2020)

| Parameter | | 2010 values | Doubling time (in years) | 10-year increase factor | value |
|---|---|---|---|---|---|
| Peak: | $C_{\text{CPU}}$ | 50.4 GF/s | 1.7 | 59.0× | 3.0 TF/s |
| | $C_{\text{GPU}}$ | 515 GF/s | | | 30 TF/s |
| Cores:[a] | $\rho_{\text{CPU}}$ | 6 | 1.87 | 40.7× | 134 |
| | $\rho_{\text{GPU}}$ | 448 | | | 18k |
| Memory bandwidth: | $\beta_{\text{CPU}}$ | 21.3 GB/s | 3.0 | 9.7× | 206 GB/s |
| | $\beta_{\text{GPU}}$ | 144 GB/s | | | 1.4 TB/s |
| Fast memory | $Z_{\text{CPU}}$ | 6 MB | 2.0 | 32.0× | 192 MB |
| | $Z_{\text{GPU}}$ | 2.7 MB[b] | | | 86.4 MB |
| Line size: | $L_{\text{CPU}}$ | 64 B | 10.2 | 2.0× | 128 B |
| | $L_{\text{GPU}}$ | 128 B | | | 256 B |
| Link bandwidth: | $\beta_{\text{link}}$ | 10 GB/s | 2.25 | 21.8× | 218 GB/s |
| Machine peak: | $R_{\text{peak}}$ | 4 PF/s | 1.0 | 1000× | 4 EF/s |
| System memory: | $E$ | 635 TB | 1.3 | 208× | 132 PB |
| Nodes ($\frac{R_{\text{peak}}}{C}$): | $P_{\text{CPU}}$ | 79,400 | 2.4 | 17.4× | 1.3M |
| | $P_{\text{GPU}}$ | 7,770 | | | 135,000 |

## Matrix-multiplication on GPU

$$\frac{p\pi}{\beta} \leq O\left(\sqrt{\frac{\gamma}{p}}\right)$$



through measurements

$a\sqrt{\dfrac{\gamma}{p}}$

Intensity

$\dfrac{p\pi}{\beta}$

Balance

Even Matmult on GPU could become memory bound!

$p\pi$

Matmult

FFT

Sparse Matvec

Approximate operational intensity

Performance (single GFLOP/s)

Balance (flop:byte)
- · 0.25
- · 0.5
- · 1
- · 2
- ● 4
- ● 8
- ● 16

**Class**
- GPU
- CPU
- Mobile
- APU  CPU+GPU
- MT  Multithreaded (old)

AMD GPU
Tesla
Cell
Fusion  x86-64
Atom
SPARC64
POWER
Blue Gene
MTA
Tegra
Niagara
iPad

Bandwidth (GB/s)

$\beta$

1. More powerful: less balance
2. Build large scale with low power processors