# A Changing Landscape

Jason Cong[2], Konstantis Daloukas[1], Nate Earnest-Noble[3], Kostas Kafousas[1], Sophia Kolak[4], Yorgos Koutsoyannopoulos[1], Bob Lucas[1], Hamed Mohamedbagherpoor[3], Francois-Henry Rouet[1], and Linghao Song[2]

[1]Ansys, [2]UCLA, [3]IBM, [4]CMU

June 22, 2023

**Ansys**

# Past era of innovation

- ## The mid-1980s saw the emergence of CMOS as a viable semiconductor technology
  - Low cost, low power, and high-volume components led to "commercial off the shelf" (COTS) strategy

- ## Unbounded demand for computing in science and national security
  - Public funding for computer science research

- ## A decade of innovation in computer architecture and parallel programming
  - Communicating sequential processes, shared memory, data parallel, systolic, dataflow, etc.

- ## Thomas Sterling ruined it for all of us
  - He invented the Linux cluster
    - Cheap hardware
    - Free software

# Déjà vu: we have entered a new era of innovation

- Dennard scaling has ended

- Moore's Law is slowing

- Yet, demand is growing

- Specialized systems are increasingly attractive
  D.E. Shaw's Anton (drug design)
  "General purpose" GPUs
  TPUs and DPUs
  Neuromorphic (IBM True North)
  Annealers for optimization (Fujitsu, Toshiba and D-Wave)
  ML start ups including:
    Cerebras (wafer scale integration)
    Samba Nova (reconfigurable data flow)
  Quantum computers

- Which of these can we exploit?
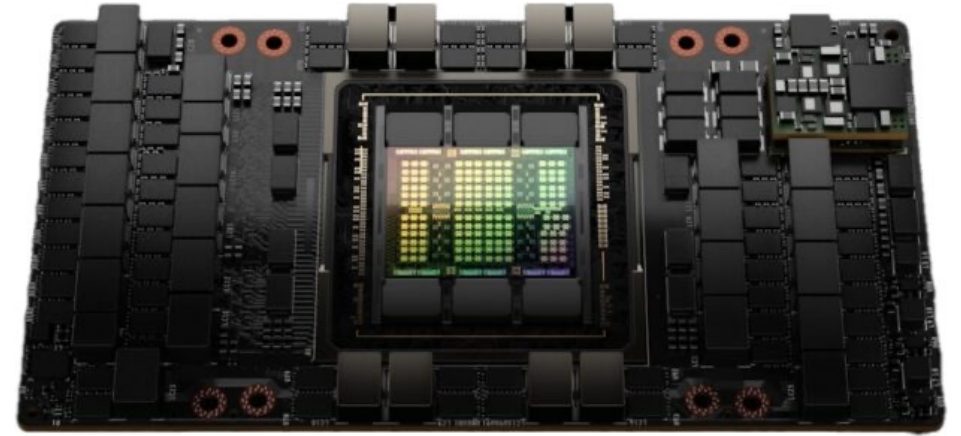


Anton



D-Wave

/Ansys

# Graphical Processing Units (GPUs)

- Familiar to this community
  - MUMPS uses them

- More computing power than their host
  - Slower clocks, but an order-of-magnitude more ALUs
  - Off-load DGEMM and other compute-bound functions
    - $O(N^2)$ data transfer versus $O(N^3)$ floating point operations
    - IBM's ESSL does this transparently
  - Write more sophisticated functions in CUDA (or HIP, or SYCL)

- More memory bandwidth than their host
  - HBM provides an order-of-magnitude more main memory bandwidth
    - An order of magnitude less memory volume
  - Off-load memory bandwidth bound functions
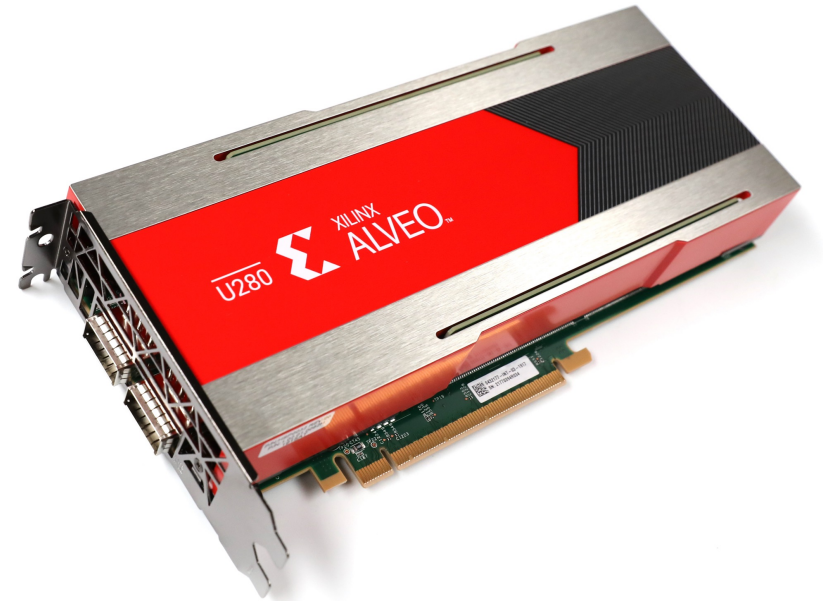    - E.g., iterative solvers dominated by sparse matrix multiply



Nvidia Hopper

**Ansys**

# Accelerated Processing Units (APUs)

- AMD devices mixing CPU and GPU
  - On the same for game consoles
  - In the same socket for HPC

- Only one physical memory
  - Upcoming MI300A CPUs and GPUs will share 128 GB of HBM
    - Will be launched at SC23
  - No DDR for the CPUs
  - Should eliminate the CPU <-> GPU data copying bottleneck

- Open questions
  - Will this enable exploiting finer-grained computations on GPUs?
  - Will we be able to use OpenMP 5 directives effectively?
  - Will kernel launch times or other overheads still limit broader utility of GPUs?

# Field Programmable Gate Array (FPGA)

- "Sea of gates"
  - First brought to market in 1985 by Xilinx
  - User programs logical units (LUTS) and their interconnect
    - Create your own custom circuit
    - Slower and more power-hungry than an ASIC
  - Widely used in networking systems
    - Reconfigure hardware in the field with software patches

- Programmable Array of Memory
  - World's first attempt to make an FPGA-based accelerator
  - DEC Paris Research Center

- SPLASH
  - First attempt in the US
  - My wife was one of four people who could program it



Xilinx U280

# How can an FPGA help?

- GPUs were initially used for compute intensive functions
  - Factor large frontal matrices

- Now GPU memory bandwidth is often more important
  - Sparse matrix – vector multiplication
  - Algebraic multigrid

- FPGAs are cheaper and have lower power consumption than GPUs

- FPGAs now have HBM too
  - Xilinx U280 has 8 GB and 460 GB/s
  - Ansys joined UCLA Prof. Jason Cong's Center for Domain-Specific Computing (CDSC)
  - Exploring preconditioned Conjugate Gradients

Ansys

# The Jacobi Preconditioned Conjugate Gradient (JPCG) Solver

**Algorithm 1** Jacobi preconditioner conjugate gradient solver for solving a linear system $\mathbf{A} \cdot \vec{x} = \vec{b}$.

**Require:**

  (1) matrix $\mathbf{A}$, (2) Jacobi preconditioner $\mathbf{M}$, (3) reference vector $\vec{b}$, (4) initial solution vector $\vec{x}_0$, (5) convergence threshold $\tau$, and (6) maximum iteration number $N_{\max}$.

**Ensure:**

  A solution vector $\vec{x}$.

1:  $\vec{r} \leftarrow \vec{b} - \mathbf{A} \cdot \vec{x}_0$
2:  $\vec{z} \leftarrow \mathbf{M}^{-1} \cdot \vec{r}$
3:  $\vec{p} \leftarrow \vec{z}$
4:  $rz \leftarrow \vec{r}^{\top} \cdot \vec{z}$
5:  $rr \leftarrow \vec{r}^{\top} \cdot \vec{r}$
6:  **for** $(0 \leq i < N_{\max}$ and $rr > \tau)$ **do**
7:    $\mathbf{a}\vec{p} \leftarrow \mathbf{A} \cdot \vec{p}$
8:    $\alpha \leftarrow rz / (\vec{p}^{\top} \cdot \mathbf{a}\vec{p})$
9:    $\vec{x} \leftarrow \vec{x} + \alpha \cdot \vec{p}$
10:   $\vec{r} \leftarrow \vec{r} - \alpha \cdot \mathbf{a}\vec{p}$
11:   $\vec{z} \leftarrow \mathbf{M}^{-1} \cdot \vec{r}$
12:   $rz\_new \leftarrow \vec{r}^{\top} \cdot \vec{z}$
13:   $\vec{p} \leftarrow \vec{z} + (rz\_new / rz) \cdot \vec{p}$
14:   $rz \leftarrow rz\_new$
15:   $rr \leftarrow \vec{r}^{\top} \cdot \vec{r}$
16: **end for**

- Solve A*x=b where A and b are known and x is unknown

- Iteratively refine errors and approach a solution

- Widely used in scientific and engineering computing, including LS-DYNA
  - Default for thermal modeling

**/\nsys**

# Callipepla

- High-Level Synthesis can take days to run
  - Therefore, a PCG accelerator cannot be uniquely designed for each sparse matrix
- Callipepla (California state bird) is streaming architecture for PCG
  - Operates on an arbitrary sparse matrix structure
- Research questions included:
  - How to support an arbitrary problem and terminate accelerated processing on the fly?
    - Main loop termination condition unknown until run time
  - How to coordinate long-vector data flow among processing modules?
    - Off-chip accesses: may waste memory bandwidth
    - On-chip accesses: which vectors and which modules can be kept ina very limited memory
  - How to reduce off-chip memory bandwidth yet maintain the double precision (FP64) accuracy?
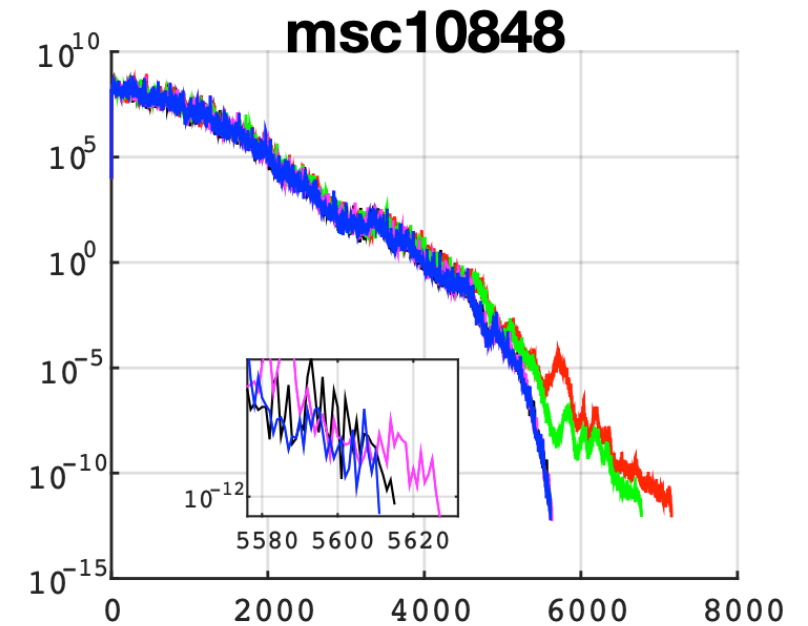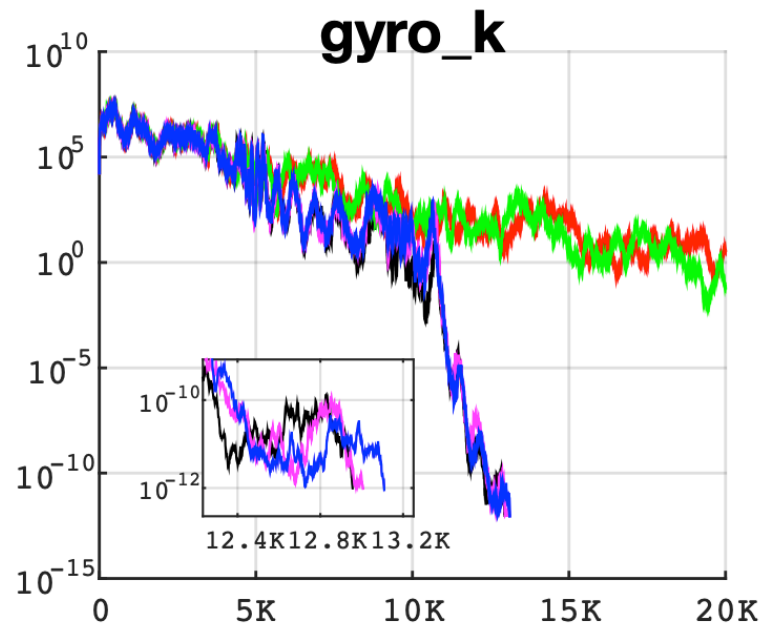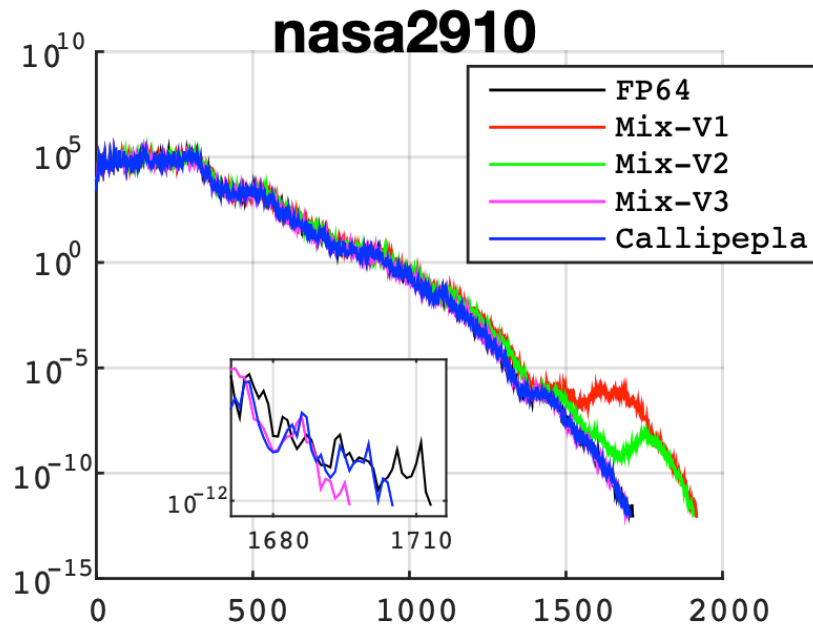    - SpMV dominates the memory bandwidth consumed

# Mixed-precision sparse matrix-vector multiplication

- Bandwidth
  - FP32 is better

- Accuracy
  - FP64 is better

THREE MIXED-PRECISION SCHEMES FOR SPMV $\vec{y} = \mathbf{A} \cdot \vec{x}$.

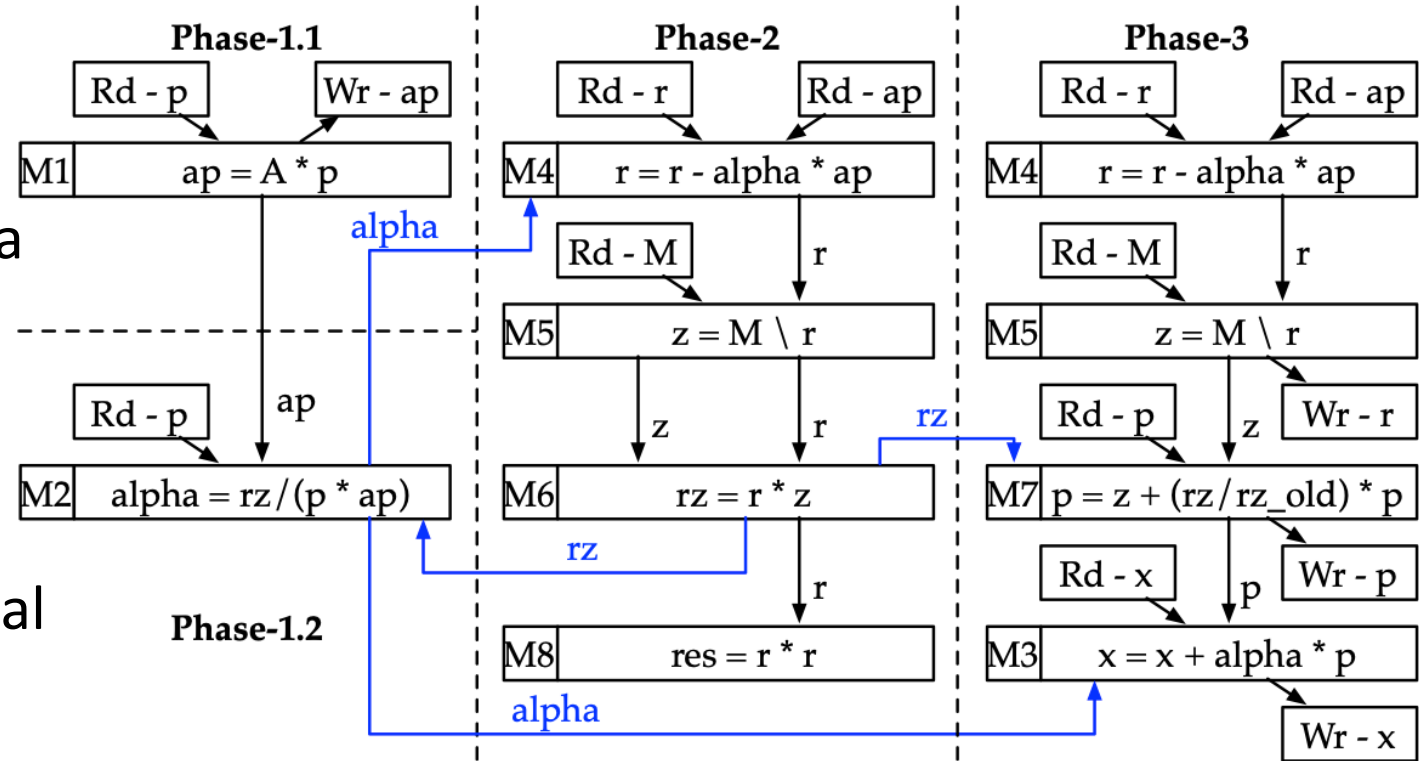|  | $\mathbf{A}$ | $\vec{x}$ | $\vec{y}$ |
|---|---|---|---|
| **Default FP64** | FP64 | FP64 | FP64 |
| **Mixed-V1** | FP32 | FP32 | FP32 |
| **Mixed-V2** | FP32 | FP32 | FP64 |
| **Mixed-V3** | FP32 | FP64 | FP64 |

©2021 ANSYS, Inc.

# Callipepalla architectural components

- U280 HBM FPGA

- Memory modules -> read/write

- Vector control modules -> control vector flows

- Computation modules -> computations, M1-M8

- SpMV: based on Serpens (DAC'22)

# Dependency & Three Computation Phases

- The computation modules -> three phases
  - A scalar dependency separates two phases
  - Each vector only rd/wr once within a phase
  - All modules share the same vectors within a phase
- Vector streaming reuse
  - To trigger the processing of individual modules
  - To overlap computation to save processing time
  - To share vectors among modules to save off-chip bandwidth

# Early results in LS-DYNA

- AWE nested cylinders benchmark
- AMD Host with Xilinx Alveo

  | Solver | 1st step | 2nd step |
  | --- | --- | --- |
  | Multifrontal | 849 | 2.1 |
  | Ref. JPCG | 408 | N/A |
  | FPGA JPCG | 43 | 50 |

- Intel Skylake (two AVX512 ALUs), first load step

  | Solver | 1 CPU | 8 MPI |
  | --- | --- | --- |
  | Multifrontal | 562 | 113 |
  | JPCG | 265 | 39 |
  | ICCG | 164 | 34 |

FPGA timing:

| | |
| --- | --- |
| Preprocessing | 28.43 |
| Downloading | 0.77 |
| Solving | 13.96 |

/Ansys

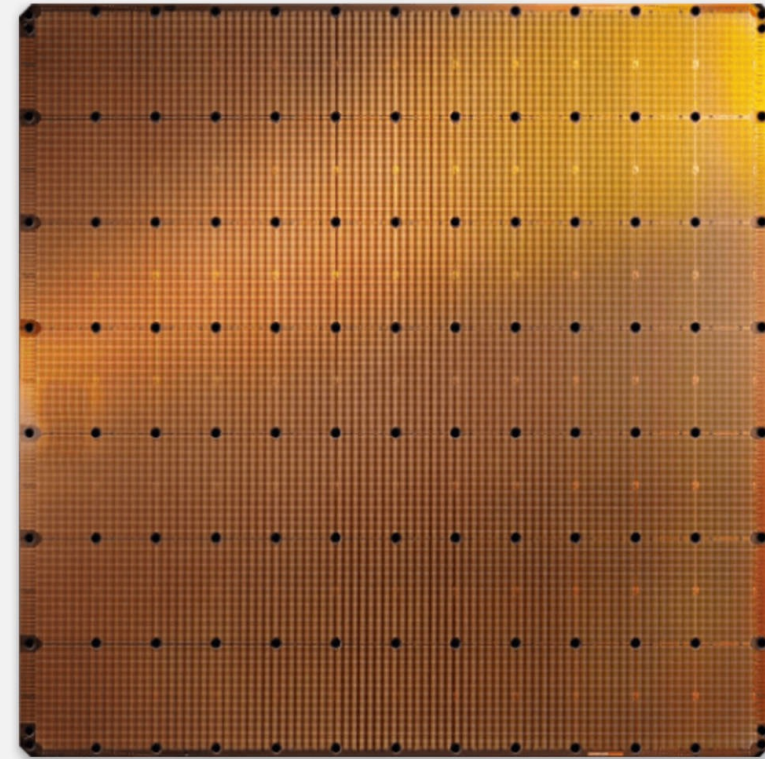# Future directions for FPGA R&D

- Multiple FPGAs per host node
  - How many PCI-e channels do you have?
  - Multiple distributed memory hosts
    - Like Nvidia GPUS, Xilinx FPGAs can talk to each other without involving their hosts

- ICCG
  - Requires triangular solves
  - Work in progress with UCLA

# Cerebras machine learning startup

- One of over a hundred hardware startups targeting training for machine learning

- Cerebras is unique in that they've gone wafer-scale!
  - Obvious good idea for decades
    - Not easy though. Gene Amdahl tried and failed
  - TSMC is their fab

- Data movement on-chip takes more energy than arithmetic.
  - 100 pJ for a floating-point operation
  - 120 pJ to move an operand 2 cm

- Data movement off-chip takes an order-of-magnitude more energy
  - 2000 pJ to move a bit to DRAM
  - The Cerebras WSE doesn't have multiple chips

**Ansys**

# Cerebras CS-2 Wafer-Scale Engine

- Good
  - Largest "chip" in the World at 46,225 mm$^2$
  - 2.6 trillion transistors
  - 850,000 32-bit cores, with ~2 PFlop/s peak
  - 20 PB/s of memory bandwidth
  - 220 Pb/s on-chip network

- Less Good
  - 40 GB of SRAM
    - Not much state at each vertex in a neural net
  - No high-level programming language
    - TensorFlow and PyTorch
  - 15 RU, 23 KW
    - More energy efficient than GPU-based systems



**Cerebras WSE-2**
46,225mm$^2$ Silicon
2.6 Trillion transistors

**Largest GPU**
826mm$^2$ Silicon
54.2 Billion transistors

# Collaboration with Cerebras and National Energy Technology Lab

- Focus on iterative solvers
  - Not enough memory for MUMPS

- Initially looking at structured, 3D grids
  - First two dimensions map easily to Cerebras processor grid
    - 3rd dimension is local memory
  - Performance of approximately 1 PFlop/s when all processors are busy

- Unstructured grids is work-in-progress

/Ansys

# Ansys and quantum computing

- Ansys is involved in quantum computing today.
  - Our Computer Aided Engineering (CAE) tools are being used to design quantum computers.
  - E.g., HFSS is used to design resonators for transmons

- Ansys needs to know if we can use quantum computers to accelerate our CAE tools.
  - We have formed a partnership with IBM to investigate quantum computing.
  - Initial focus on graph partitioning.

- Longer term, Ansys expects quantum computers will enable new CAE markets that are inconceivable today.
  - Simulation of quantum dynamics.

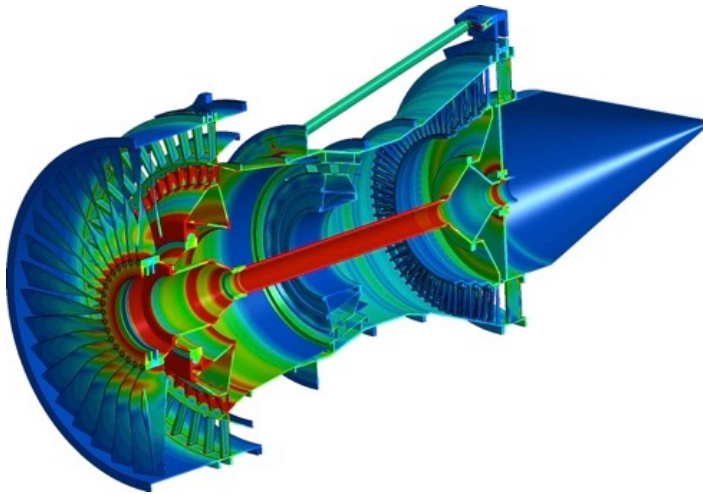# The relative importance of sparse matrix reordering

We reorder sparse matrices to reduce the cost of factorization
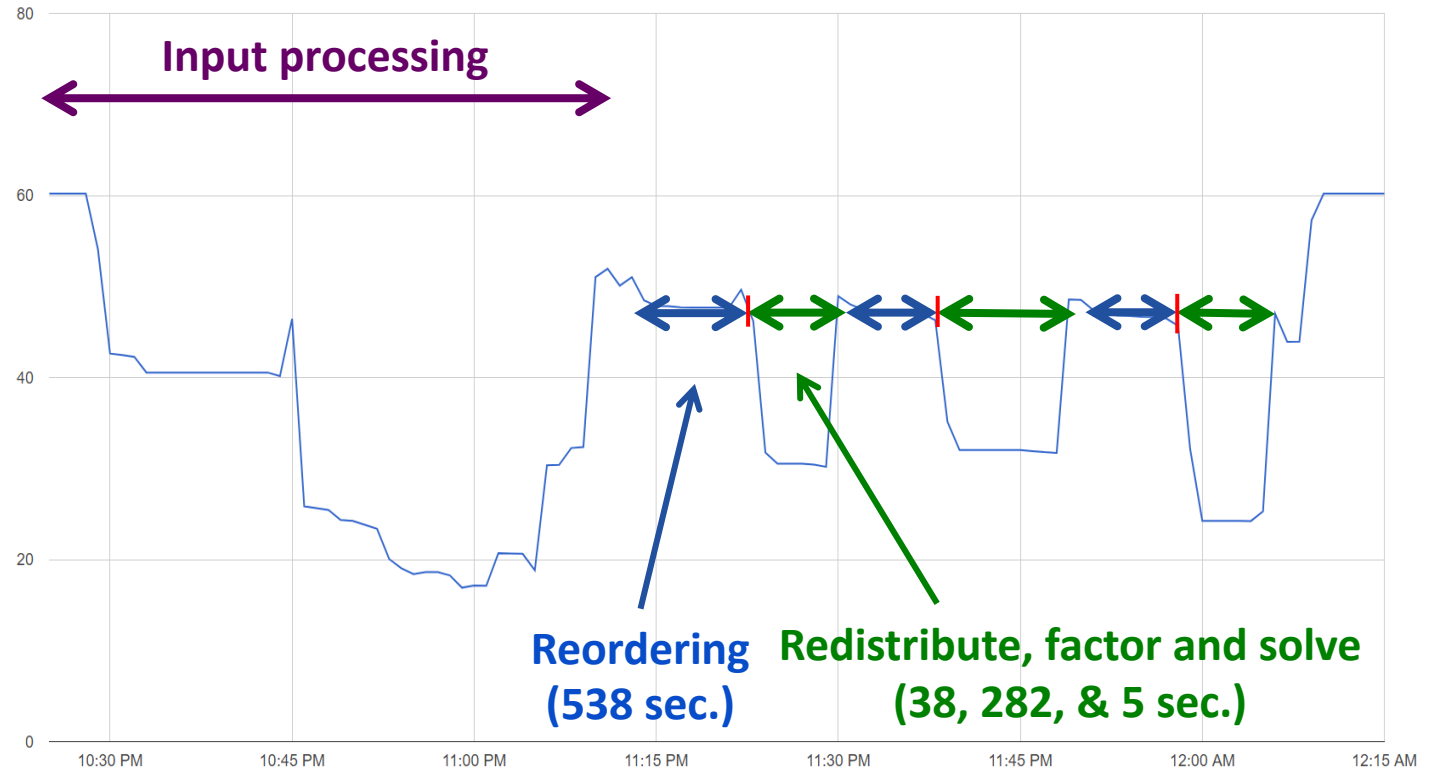
- Minimizing cost is NP-complete
- Approximate with nested dissection

As you add processors, eventually reordering and factorization dominate the run time.

- Factorization scales better



Rolls-Royce Representative Engine Model (REM)

Three simulated load steps of the Rolls-Royce REM.



Input processing

Reordering (538 sec.)

Redistribute, factor and solve (38, 282, & 5 sec.)

LS-DYNA memory usage (rank 0) vs. time on NCSA Blue Waters
256 MPI ranks, 8 threads each

**Ansys**

# Reordering with nested dissection

## LS-GPart nested dissection algorithm

For each domain:

    Coarsen graph

    Partition coarsened graph

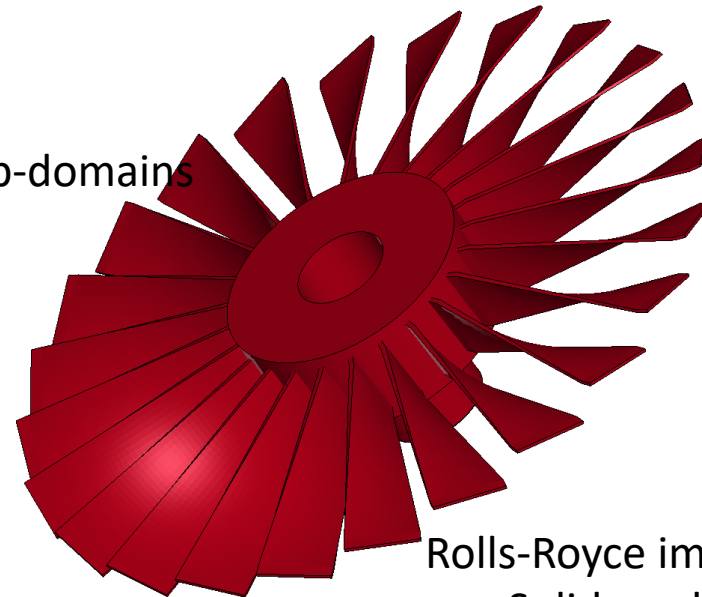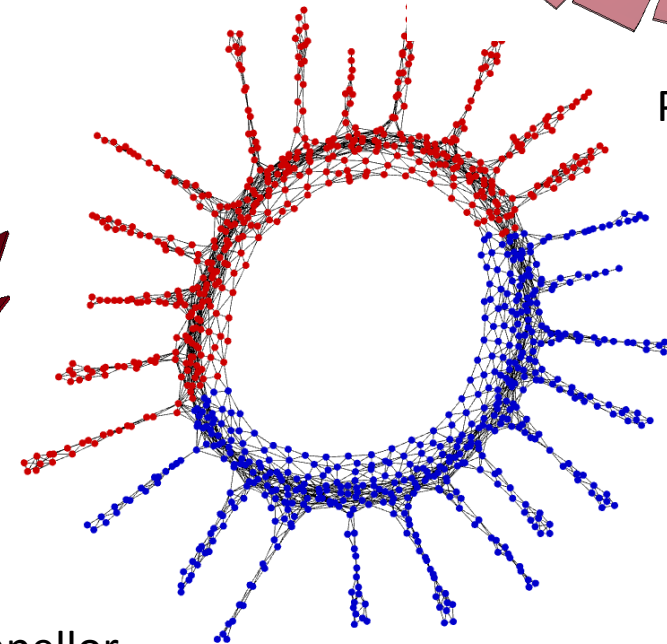        NP-complete

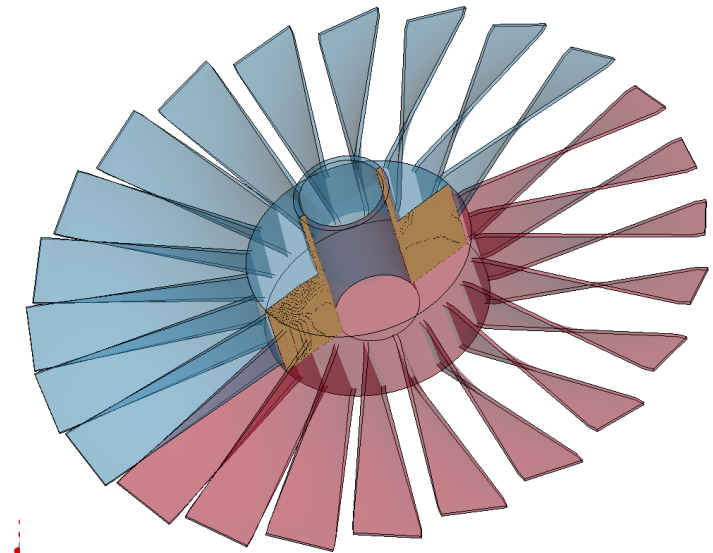    Project to fine graph

    Refine

    Remove separator

    Recurse on resulting sub-domains
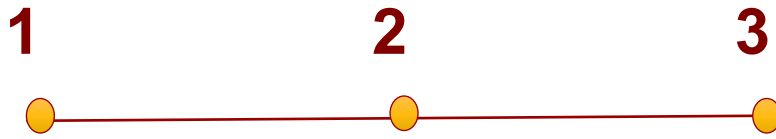
Partitioned model

Rolls-Royce impellor
Solid model

Coarse graph

# Graph partitioning as an Ising model

$$H = A \left( \sum_{i \in V} \sigma_i \right)^2 + B \sum_{ij \in E} \frac{1 - \sigma_i \sigma_j}{2}$$

$\underbrace{\phantom{\left( \sum_{i \in V} \sigma_i \right)^2}}_{\substack{\text{penalize config.} \\ \text{that doesn't bisect}}}$ $\underbrace{\phantom{\sum_{ij \in E} \frac{1 - \sigma_i \sigma_j}{2}}}_{\substack{\text{penalty } B \text{ for} \\ \text{each edge btwn} \\ \text{two subsets}}}$

**1**          **2**          **3**

Solutions$^{\mathsf{T}}$          **H**          All possible solutions          Energy

|       |   |   |   |   |   |   |   |
|-------|---|---|---|---|---|---|---|
| - - - |   |   |   |   |   |   |   |

$\begin{matrix} - & - & - \\ - & - & + \\ - & + & - \\ - & + & + \\ + & - & - \\ + & - & + \\ + & + & - \\ + & + & + \end{matrix}$ * $\begin{matrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \end{matrix}$ * $\begin{matrix} - & - & - & - & + & + & + & + \\ - & - & + & + & - & - & + & + \\ - & + & - & + & - & + & - & + \end{matrix}$ = $\begin{matrix} 5 \\ 1 \\ 5 \\ 1 \\ 1 \\ 5 \\ 1 \\ 5 \end{matrix}$
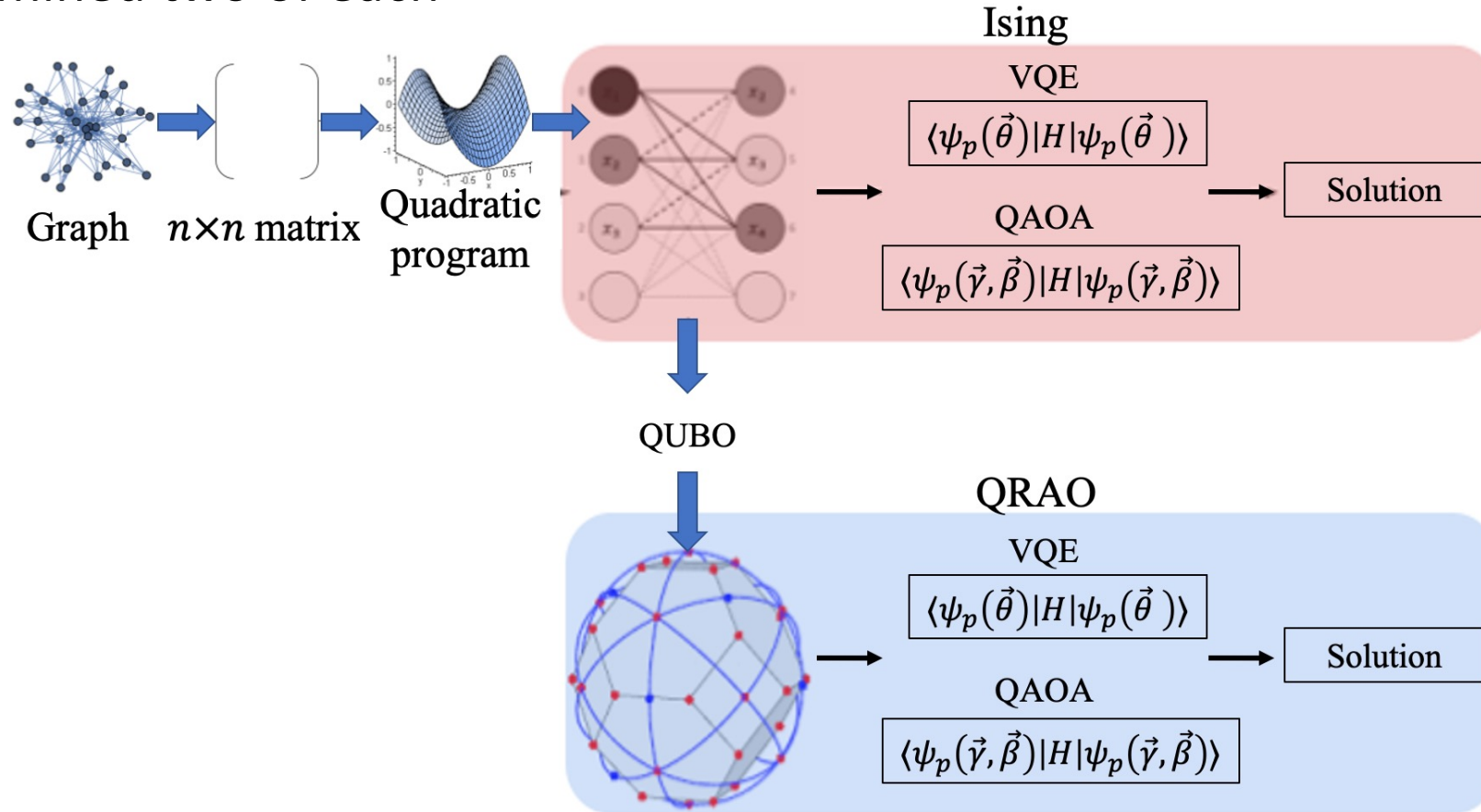
The objective is to minimize the product of

$x^{\mathsf{T}} * H * x$

Program H so as to add energy penalties to configurations you don't want
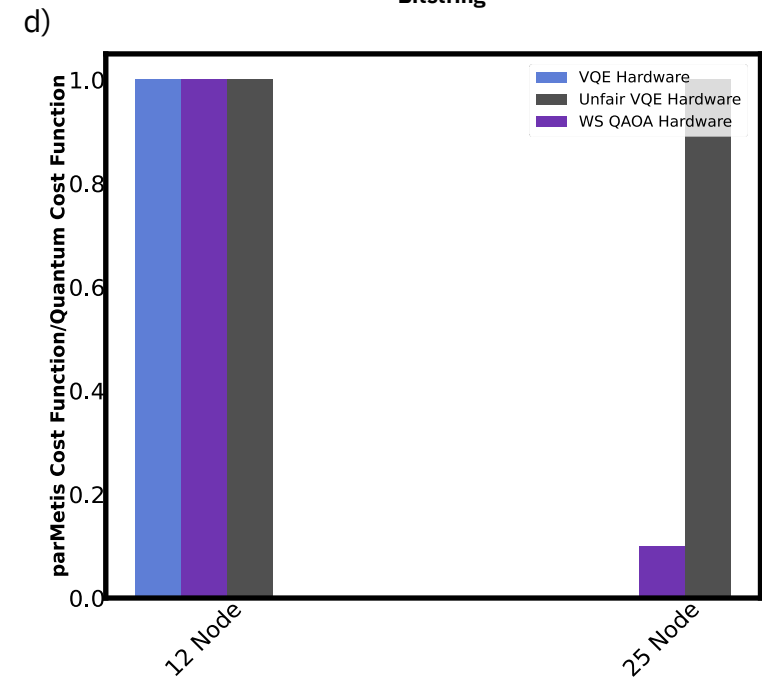
J Phy A19 1605 (1986)
Frontiers in Physics 2, 5 (2014)

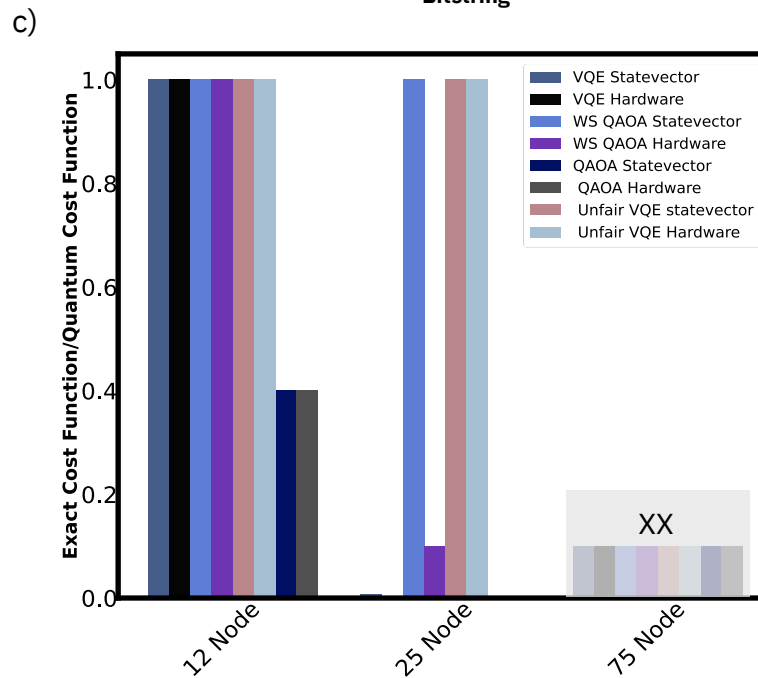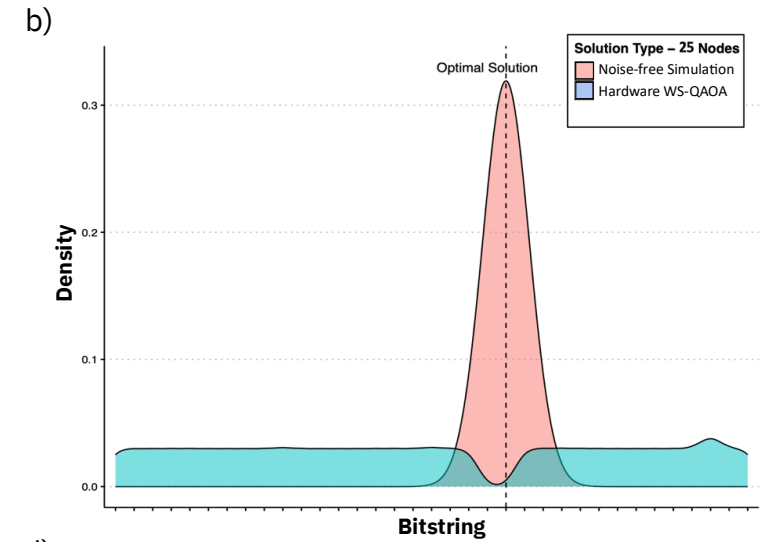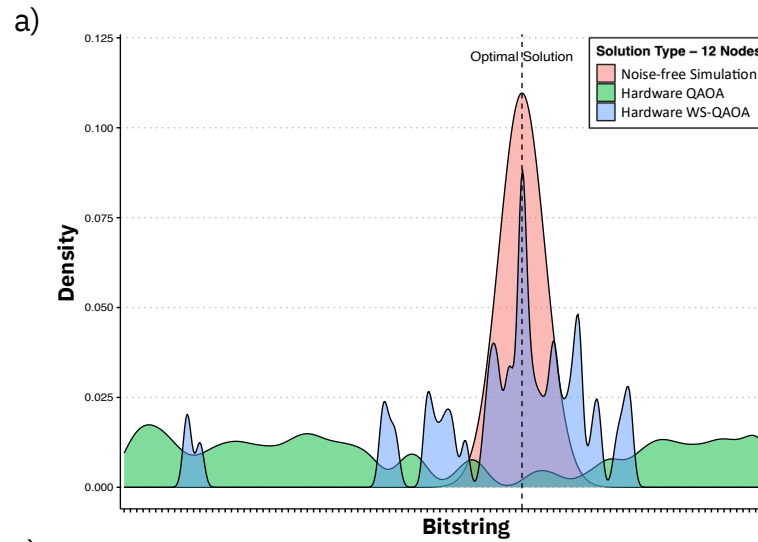# Encoding and solutions considered

- The are multiple possible encodings and multiple possible solution algorithms
- We examined two of each



©2021 ANSYS, Inc.

# Early results

- Results on 12 and 25-qubit IBM quantum computers

# Summary

Machine learning has increased the demand for computing power
- Demand is growing faster than Moore's Law
- Meanwhile, Moore's Law itself it slowing down

We have entered a new era of innovation in computing technology
- Primarily privately funded this time around
- Driven by advertising

This presents us with a rich new set of potential accelerators to explore
- Look beyond GPUs