# HPC Challenges in Plasma Physics

### **Frank Jenko**

Max Planck Institute for Plasma Physics, Garching Technical University of Munich University of Texas at Austin

**IHPCSS 2019 – Kobe, Japan – July 9, 2019** 

### Global electricity needs will keep increasing

Energy Modeling Forum 22 100 models from 15 research groups (Clarke 2009)



### Fusion energy in the laboratory



Still, temperatures of about 100 million degrees are required! Thus, we are dealing with a fully ionized gas (plasma).

### Magnetic confinement of fusion plasmas

Charged particles basically follow magnetic field lines

Helically twisted field lines span nested magnetic surfaces

Such an axisymmetric device is called "tokamak"



### Fusion research: Towards burning plasmas



### The international ITER project



### Goal: 500 MW of fusion power www.iter.org **ITER PROJECT: International Cooperation Russian Federation** Europea Union USA Korea China Japan ITER Partners The seven parties involved in the ITER construction represent more than 50 % of the world's population

### **ITER construction site in Southern France**



### From trial-and-error to predict-first: The role of High Performance Computing

## At the forefront of supercomputing since the 70's

### **NERSC HISTORY**

### **Powering Scientific Discovery Since 1974**

Contact: Jon Bashor, jbashor@lbl.gov, +1 510 486 5849

The oil crisis of 1973 did more than create long lines at the gas pumps - it jumpstarted a supercomputing revolution.

The quest for alternative energy sources led to increased funding for the Department of Energy's Magnetic Fusion Energy program, and simulating the behavior of plasma in a fusion reactor required a computer center dedicated to this purpose. Founded in 1974 at Lawrence Livermore National Laboratory, the Controlled Thermonuclear Research Computer Center was the first unclassified supercomputer center and was the model for those that followed.

Over the years the center's name was changed to the National Magnetic Fusion Energy Computer Center and later the National Energy Research Supercomputer Center (NERSC). In 1983 NERSC's role was expanded beyond the fusion program, and it becan providing general computing services to all of the programs funded by the DOE Office of Energy Research (now the Office of Science). The current name was adopted in 1996 when NERSC relocated to Lawrence Berkeley National Laboratory and merged with Berkeley Lab's Computing Sciences program. The name change — from "Supercomputer Center" to "Scientific Computing Center" — signaled a new philosophy, one of making scientific computing more productive, not just providing supercomputer cycles.



### Towards a virtual fusion plasma

Increasing fidelity & modeling capability with increasing computing power



#### Gigaflops

Core: ion-scale electrostatic physics in simplified geometry



#### Beyond

Whole device modeling of all relevant fusion science

Acknowledgements: ECP

Goals: prepare and interpret **ITER** discharges, guide the development of **power plants** 

### The multiscale, multiphysics challenge



Many nonlinear interactions; we cannot use a simple "superposition principle"

## A multi-fidelity approach

An example:





- High-fidelity models provide reliable predictive capability
- Lower-fidelity models foster high-throughput computing
- Both are needed together

### A high-fidelity model for determining turbulent transport (i.e., the energy confinement time): The GENE code

### Fluid models don't work – use (gyro-)kinetics!

Hot and/or dilute plasmas are only weakly collisional: 6D Vlasov-Maxwell equations

$$\frac{\partial f_{\alpha}}{\partial t} + \mathbf{v} \cdot \nabla f_{\alpha} + \frac{q_{\alpha}}{m_{\alpha}} \Big[ \mathbf{E} + \frac{\mathbf{v} \times \mathbf{B}}{c} \Big] \cdot \nabla_{v} f_{\alpha} = 0 \quad \alpha = \text{particle species}$$

... from the Liouville equation via the BBGKY hierarchy



 $f_{\alpha} = f_{\alpha}(\mathbf{x}, \mathbf{v}, t)$ 

#### Strong background magnetic field:

Eliminate fast gyromotion; consider slow dynamics of **guiding centers** 

$$f = f(\mathbf{X}, v_{\parallel}, \mu; t)$$

$$\frac{\partial f}{\partial t} + \dot{\mathbf{X}} \cdot \frac{\partial f}{\partial \mathbf{X}} + \dot{v}_{\parallel} \frac{\partial f}{\partial v_{\parallel}} = 0$$

Reduction of effort by ~12 orders of magnitude (elimination of irrelevant spatio-temporal scales & reduction from 6D to 5D)

Additional gain from using field-aligned coordinates

## The gyrokinetic code GENE



- First GENE publication: Jenko et al., Physics of Plasmas 2000 (>600 citations WoS)
- More than 100,000 lines of source code, plus 200,000 lines for pre-/post-processing
- Based on numerical methods from Computational Fluid Dynamics
- Open source policy
- World-wide user base: genecode.org support@genecode.org

Part of an ecosystem of codes for fusion research

Also used in astrophysics



Global Gyrokinetic Simulation of Turbulence in ASDEX Upgrade



gene.rzg.mpg.de

### Some background on GENE

- The underlying nonlinear PDEs are discretized on a fixed grid in 5D phase space
- Apply CFD-type (mix of spectral, finite difference, and finite volume) methods
- Explicit time-stepping facilitates scalability
- Time step is maximized during initialization



Runge–Kutta– Chebychev schemes

Doerk & Jenko CPC 2014

- 5D domain decomposition
- Primarily pure MPI or hybrid MPI/CUDA/OpenACC programming model
- Auto-tuning (optimal subroutines and processor layout)

### Verification and validation



### Science highlights (just two examples)

Prediction of relevant contributions to turbulent transport at very small – hiterto neglected – spatio-temporal scales (experimentally confirmed)

Indications that the physics of the "pedestal region" in ITER may differ from that in present-day devices





### Strong scaling of GENE on Titan



Tilman Dannert



Optimizing large-scale codes (like GENE) on pre-exascale systems may take person-years

### From Titan to Summit (2018)

2,282,544 cores 187 PF (peak)

**Rank System** 

IBM's Power9 processors NVIDIA's Volta GV100 GPUs

#### Much faster data motion



1 **Summit** - IBM Power System AC922, IBM POWER9 2,282,544 22C 3.07GHz, NVIDIA Volta GV100, Dual-rail Mellanox EDR Infiniband (/system/179397), IBM DOE/SC/Oak Ridge National Laboratory (/site /48553) United States

Most powerful supercomputer in the world (Top 500 List, June 2018)

122,300.0

187,659.3 8,806

### Some software engineering aspects

Version control system

Repository hosted by



Web interface installed at MPCDF – contains wiki, CI/CD, interface to bug tracker

Development repository managed via



https://gitlab.com

**Regression tests** are performed automatically once per day/week

User repository at https://gitta.rzg.mpg.de/GENE

release tags and master branch mirrored from development repository

Source code documentation via Doxygen and FORD (FORtran Documenter)

### Beyond brute force: Sparse grid combination technique

EXAHD - An Exa-Scalable Two-Level Sparse Grid Approach for Higher-Dimensional Problems in Plasma Physics and Beyond



## Sparse grid combination technique

#### **Cartesian grid**

- Regular data structure
- Huge number of grid points in high dimensions "curse of dimensionality"

### Sparse grid combination technique

- Good approximation of the Cartesian grid solution
- Smaller number of grid points
- Existing code (GENE) can be used more or less as it is

Resolution: 33 grid points per dimension	2D	5D
Cartesian grid	1,089	39,135,393
Combination tech.	641	206,358



### Combining GENE eigenvalue simulations



- Combination of eigenmodes or eigenvalues
- Here: 2D projections (vp-z) of eigenmodes



With C. Kowitz & M. Hegland

### A new level of parallelism

#### **Two-level** parallelism

- Massively parallel GENE runs for independent grid setups from the combination technique
- Run times of the instances tend to vary strongly

#### **Optimize the load balance**

- A simple **load model** estimates the runtime required for each grid setup
- A **scheduler** creates an optimal load balancing to minimize idle cores





Node 1 Node 2 Node 3



### Spin-off: Algorithmic fault tolerance

#### Hardware failures (on many Mcores)

- Standard: The whole simulation has to be restarted from the last checkpoint file
- In the combination technique, only a single GENE instance would crash

#### Two ways to handle the failure

- The combination technique recovers an approximation
- Only a single GENE instance is rerun

#### Such techniques may be very useful on emerging exascale architectures



Lossy compression of scientific data in GENE simulations of plasma turbulence

# Data outgrows compute, calls for data reduction techniques



Lossy compression enables greater reduction, but it is often met with skepticism by scientists

P. Lindstrom, IPAM Workshop, UCLA (10/15/18)

- Large improvements in compression are possible by allowing even small errors
  - Simulation often computes on meaningless bits
    - Round-off, truncation, iteration, model errors abound
    - Last few floating-point bits are effectively random noise
- Still, lossy compression often makes scientists nervous
  - Even though lossy data reduction is ubiquitous
    - **Decimation** in space and/or time (e.g. store every 100 time steps)
    - Averaging (hourly vs. daily vs. monthly averages)
    - Truncation to single precision (e.g. for history files)
  - State-of-the-art compressors support error tolerances





### Turbulence provides an interesting test case

- Turbulent dynamics is inherently nonlinear and high-dimensional
- It is characterized by a mix of disorder and order
- The detailed dynamics is chaotic, but its statistical properties tend to be robust
- It may be possible to trade accuracy for efficiency on the statistical level



### GENE-ZFP compression numerical experiment

- Emulate simulations on compressed arrays of grid-based data
- Steps in time loop:
  - 1. Compress and decompress 5D distribution function
  - 2. Compute time step
- Grid: *nx0×ny0×nz0×nv0×nw0 = 60×128×8×32×10*
- Tested ZFP modes (4D [*nx0×ny0×nz0×nv0*] compression):
  - 1. Fixed rate (fixed number of bits per floating point value)
  - 2. Fixed accuracy (fixed absolute error tolerance)
  - 3. Fixed precision (fixed number of uncompressed bits per value)

### Heat flux as figure of merit



Heat flux in normalized units, averaged over space, as a function of time

Shown is a quasi-stationary saturated turbulent state

Plots for the reference and fixed rate (10 bits per double) compression runs

# Heat flux for different compression modes and compression ratios



Time-space averaged heat fluxes, with standard deviations:

- Fixed rate scan (bits per double): (unstable for 4) 6, 8, 10, 20, 25, 30, 35, 40, 50, 60
- Fixed accuracy scan (absolute error tolerance): (wrong results for 1e-2) 1e-3 ... 1e-10, 1e-12, 1e-14, 1e-16
- Fixed precision scan (number of uncompressed bits per value): (unstable for 10) 15, 20, 25, 30, 35, 40, 50, 60

### Lossy compression in GENE simulations of plasma turbulence Denis Jarema, Peter Lindstrom & Frank Jenko

- **Goal:** Explore potential for data reduction in the GENE simulations; would help to reduce data motion and thus to increase code performance
- **Incentive:** We care mainly about *statistical* properties in a quasi-stationary saturated turbulent state, which may be quite robust (here: heat flux)
- **Results:** Compression ratios *up to ~10* with acceptable loss of quality
- **Conclusion:** Study points to *enormous potential* for data reduction in GENE

An ambitious project: Coupling two high-fidelity codes to create a whole device model

### 1<sup>st</sup> U.S. exascale system (2021)

#### ALCF 2021 EXASCALE SUPERCOMPUTER – A21



### U.S. DOE Exascale Computing Project (ECP)

#### A holistic approach

Application Development	Software Technology	Hardware Technology	Exascale Systems
Science and mission applications	Scalable software stack	Hardware technology elements	Integrated exascale supercomputers
	Correctness Visualization Data Analysis Applications Co-Design Programming models, development environment, and numitrees System Software, resource management threading standarding, motioning, and control Node OS, nurtimes Headware interface		

The ECP is a 7-year project with a cost range of \$3.5B – \$5.7B

### Fusion Energy Application Development (2016-)

### **10-year Goal:** A First-Principles-Based Whole Device Model that Covers the Full Space/Time Scales of a Reactor

- XGC full-f particle-in-cell technique with continuity across separatrix
- GENE continuum delta-f capability for core



### 1<sup>st</sup> step: Code benchmarking

### Cross-verification between GENE, XGC, and ORB5: linear ITG instability (S. Ku, G. Merlo, E. Lanti (SPC, EPFL))



Codes agree within 10% for all modes considered





### 3<sup>rd</sup> step: GENE-XGC coupling







Combining computation with data analytics and machine learning for plasma physics

### An important, timely topic of broad interest

"Science at extreme scales: Where big data meets large-scale computing"



Interdisciplinary Long Program @UCLA September 12 - December 14, 2018 200+ participants, 50+ long-term participants

#### Speaker list includes:

- Yann LeCun (Director of AI Research @Facebook)
- Emmanuel Candes (Stanford University)
- Rajat Monga (Google)
- Matthias Troyer (Microsoft)
- James Sexton (IBM)
- Adrian Tate (Cray)
- Alan Lee (AMD)

### Transformative Enabling Capabilities for fusion

Advanced Algorithms – Advanced algorithms will transform our vision of feedback control for a power-producing fusion reactor. The vision will change from one of basic feasibility to the creation of intelligent systems, and perhaps even enabling operation at optimized operating points whose achievement and sustainment are impossible without highperformance feedback control. The area of advanced algorithms includes the related fields of mathematical control, machine learning, artificial intelligence, integrated data analysis, and other algorithm-based R&D. Given the pace of advances, control solutions that establish fusion reactor operation will become within reach, as will the discovery and refinement of physics principles embedded within the data from present experiments. This TEC offers tools and methods to support and accelerate the pace of physics understanding, leveraging both experimental and theoretical efforts. These tools are synergistic with advances in exascale and other high-performance computing capabilities that will enable improved physics understanding. Machine learning and mathematical control can also help to bridge gaps in knowledge when these exist, for example to enable effective control of fusion plasmas with imperfect understanding of the plasma state.

### Deep Learning for real-time plasma control

Plasma tomography: Use CNNs to reconstruct cross-section from projections



Traditional inversion schemes:

- Varying runtime
- Dependence on additional data (magnetic equilibrium)
- Not real-time capable

Deep learning for plasma tomography using the bolometer system at JET PhD student at IPP

Francisco A. Matos<sup>a</sup>, Diogo R. Ferreira<sup>a,\*</sup>, Pedro J. Carvalho<sup>b</sup>, JET Contributors<sup>1</sup>

### Unsupervised inversion with Deep Learning



Encode (weak) prior knowledge directly in the NN architecture (in lieu of minimizing an additional regularization term):

- Positivity
- Locality (smoothness)



#### Data and Backprojection (a.u.)



Similarity to financial data analysis, earthquake prediction etc. DL for real-time disruption prediction



J. Kates-Harbeck et al., Nature 2019

### Integrating scientific knowledge into ANNs

To maximize success, *Scientific ML* should not be physics-agnostic

- Physics-guided design of ANNs
- Physics-guided learning of ANNs
- Combining data- and physics-based models

#### Potential to improve accuracy, efficiency, interpretability, generalizability



### Innovative ideas (Nils Thuerey, TUM): Accelerating fluid simulations with Deep Learning



### Conclusions



Overarching goal: Contribute to the gradual development of a validated predictive capability ("virtual fusion plasma"), helping to accelerate fusion energy research

A beautiful example of how fascinating science on some of the world's largest supercomputers contributes to solving grand challenges facing society

The development and application of GENE illustrate the **fascinating challenges and opportunities at the interface of applied mathematics, computer science & physics**