



Parallel in time and object storage with CORTX

Debasmita Samaddar

CORTX Architect series Talk 6 May 2021

- Introduction – Fusion and UKAEA – who we are?
- Parallel in Time schemes
- Exascale, Tiered storage and CORTX

Take home message:

Data is a big challenge in Plasma fusion and tiered object storage such as CORTX can provide a solution.

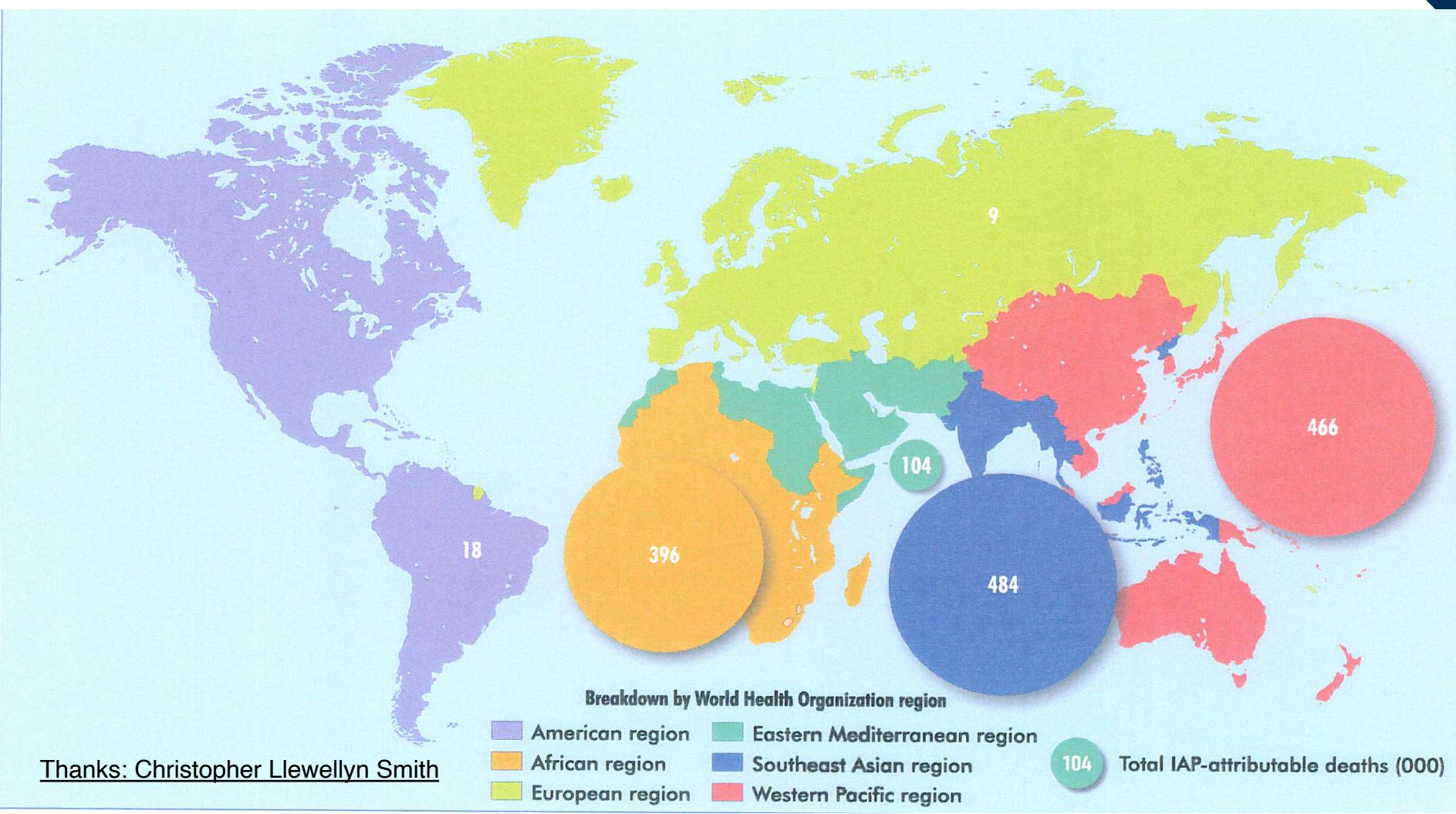
Energy requirement

- 1) **The world uses a lot of energy** – average power consumption is 2,350 Watts per person
(world energy [electricity] market ~ \$4.5 trillion [\$1.5 trillion] pa)
- 2) **World energy use is expected to grow 50% by 2030**
 - growth necessary in developing countries to lift billions of people out of poverty
- 3) **80% is generated by burning fossil fuels**
 - *climate change & debilitating pollution*
 - *which won't last for ever*

Need major new sources of clean energy - this will require fiscal measures, regulation and new technology

Thanks: Christopher Llewellyn Smith

Deaths per year (1000s) caused by indoor air pollution (biomass 85% + coal 15%); total is 1.5 million – over half children under five



Thanks: Christopher Llewellyn Smith

Energy sources

■ World's primary energy supply (approximate):

- 80 % - burning fossil fuels (44% oil, 31% coal, 26% natural gas)
- 10% - burning combustible renewables and waste
- 5% - nuclear
- 5% - hydro
- 0.5% - geothermal, solar, wind, . . .

Fossil fuels:

* generating debilitating pollution

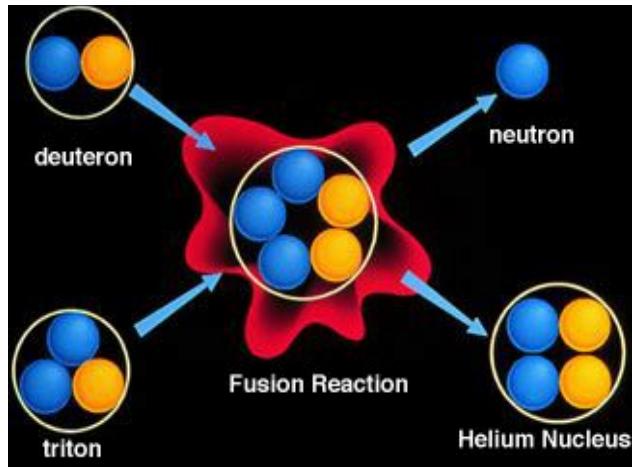
(300,000 coal pollution deaths pa in China (World Bank estimated in 1995)

* and driving potentially catastrophic climate change

* and will run out sooner or later

Thanks: Christopher Llewellyn Smith

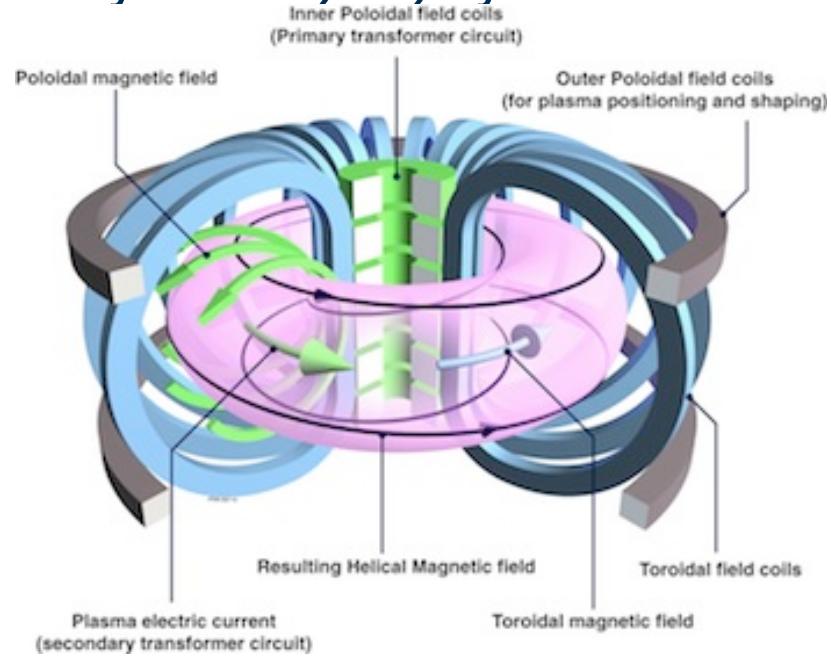
Fusion: source of energy



Fusion: power generation in the stars.

Possible in laboratories? Yes

Possible on grids for everyday life on earth?

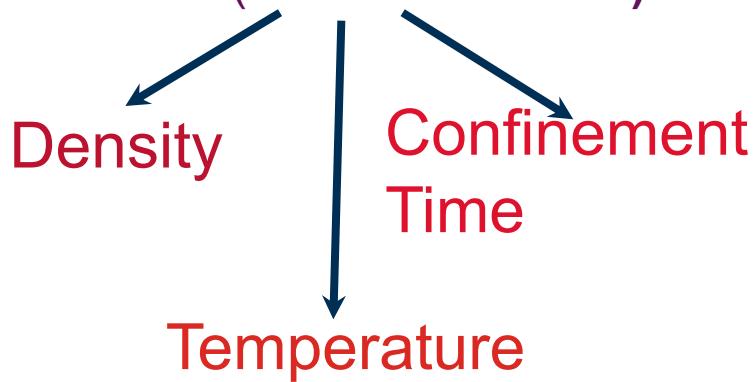


Fusion:

- a) Reaction requires temperatures $\sim 10^8$ K,
- b) Virtually unlimited – Deuterium abundant in sea water. Tritium produced from lithium which is also abundant in nature.
- c) Held together by magnetism.
- d) Relatively environmentally friendly : Waste with radioactive toxicity < 100 years

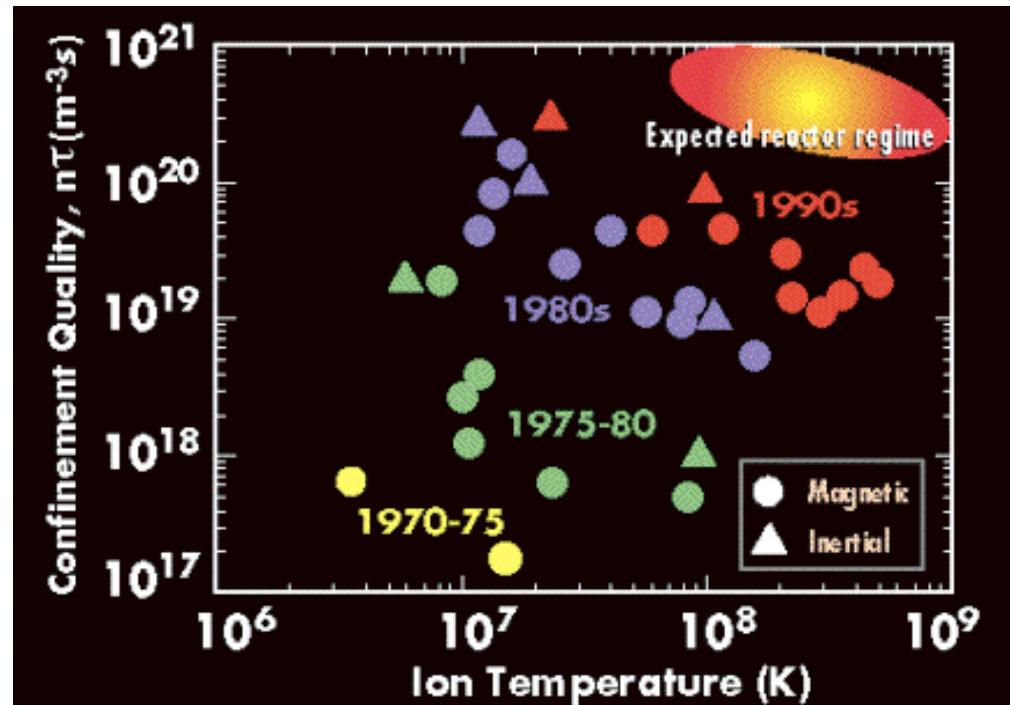
Fusion

- Typically, plasma at 100 million K
- Maximize $nT\tau$ (Lawson criterion)

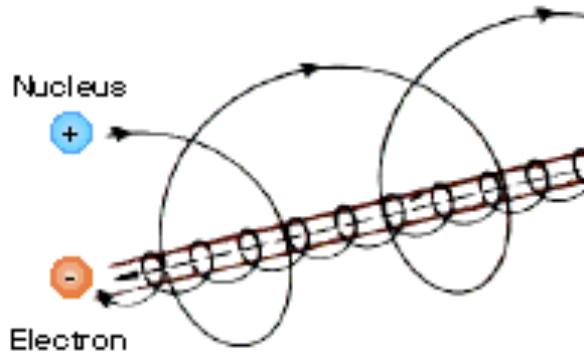


Confinement :

- Gravitational confinement - as in stars
- Magnetic Confinement of plasma
- Inertial confinement

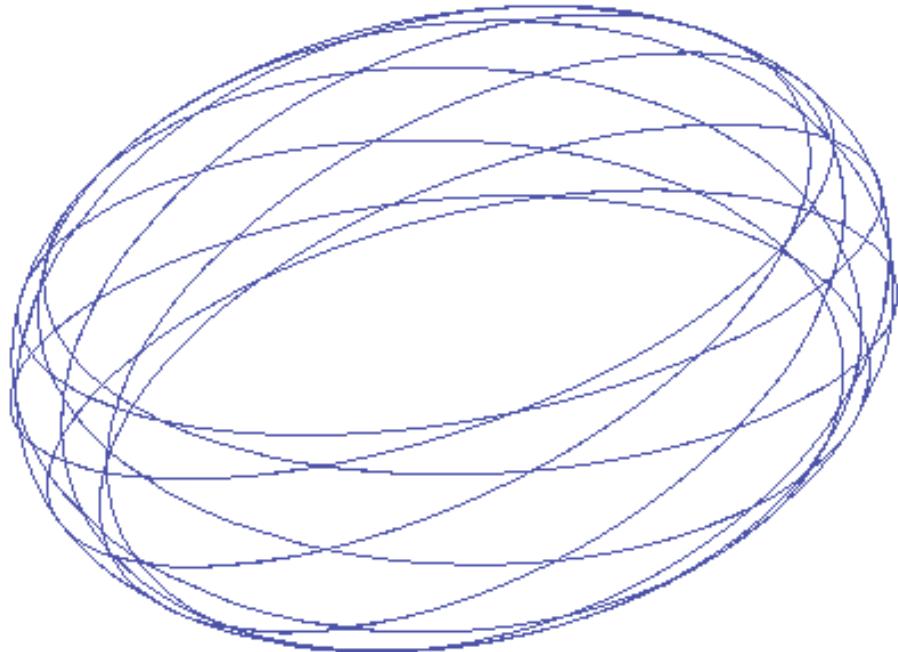


Magnetic Confinement



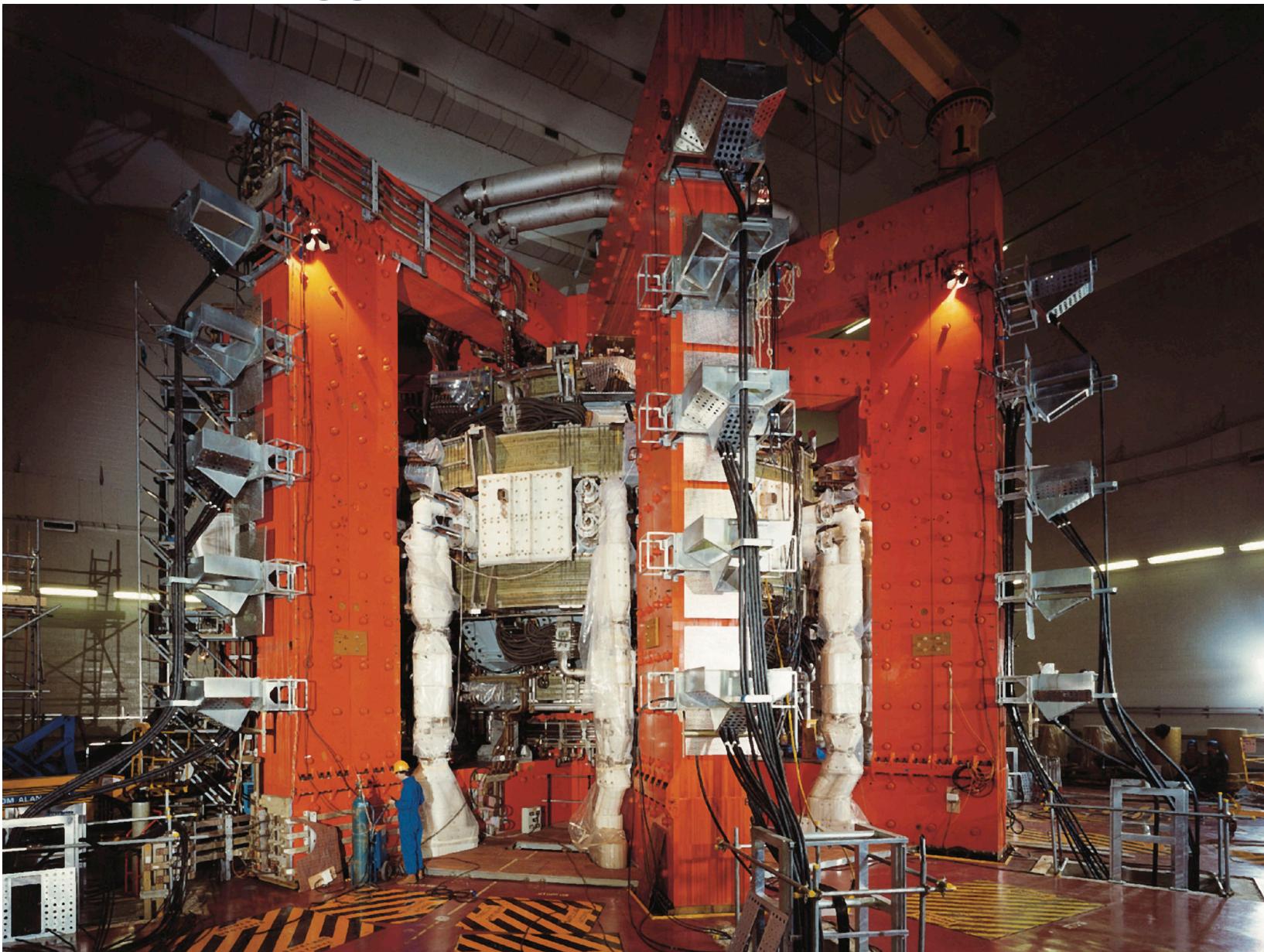
<http://ffden-2.phys.uaf.edu/>

Most devices are toroidal shaped.

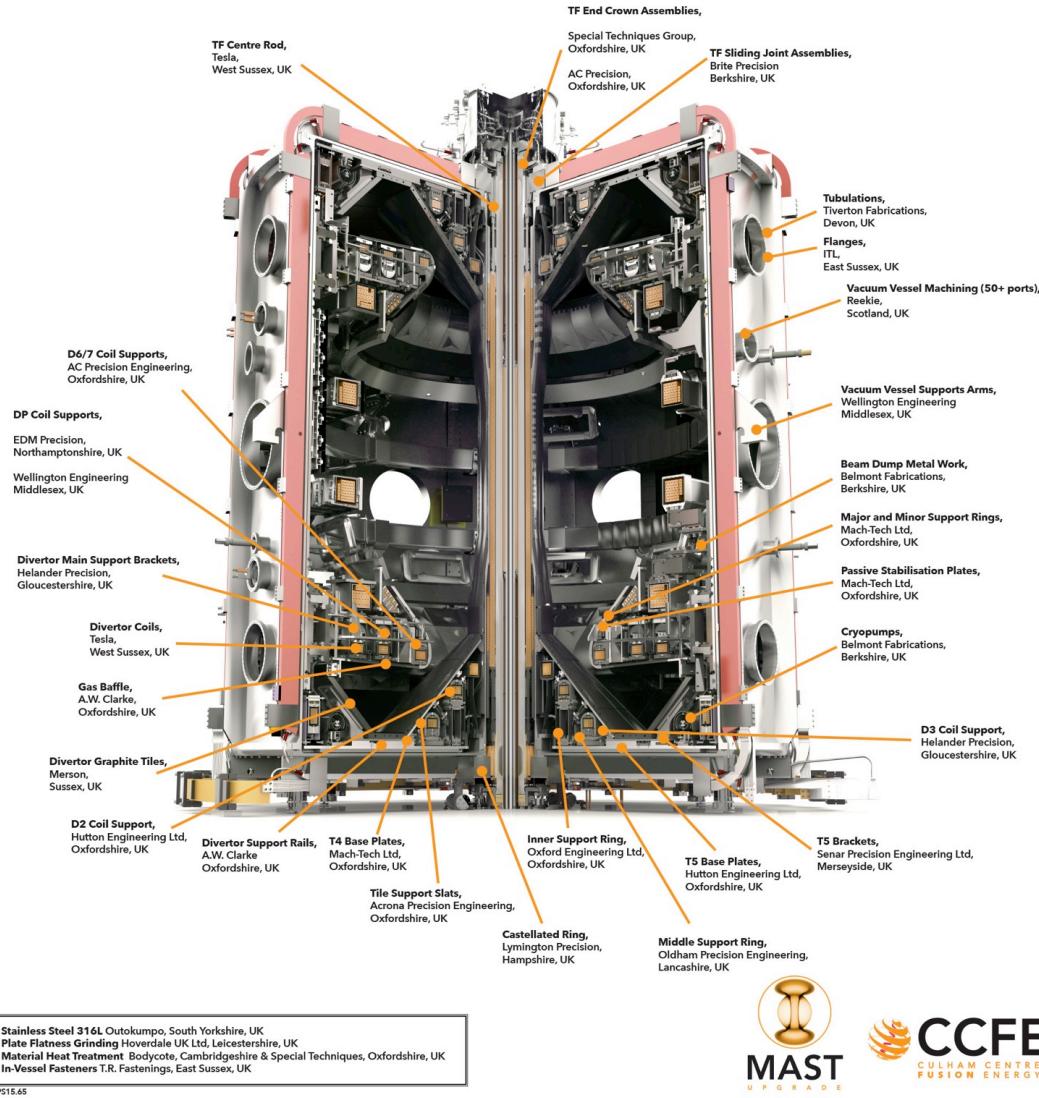
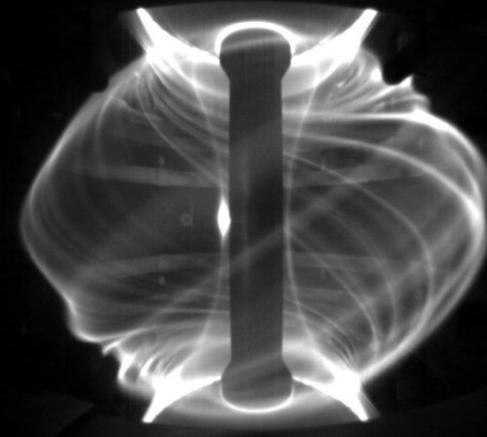


Used in ITER, DIII-D, JET, JT-60 ...

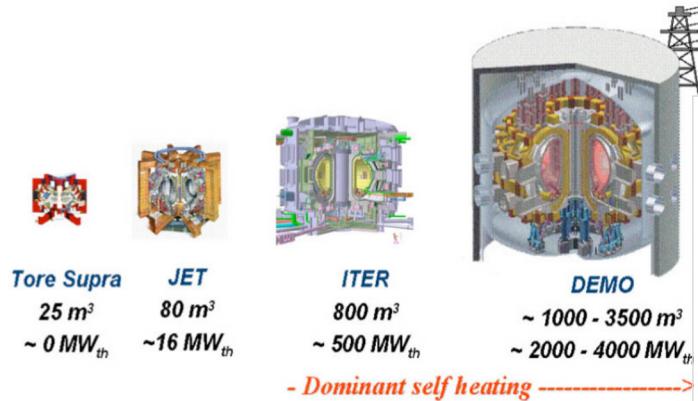
JET – biggest TOKAMAK in the world



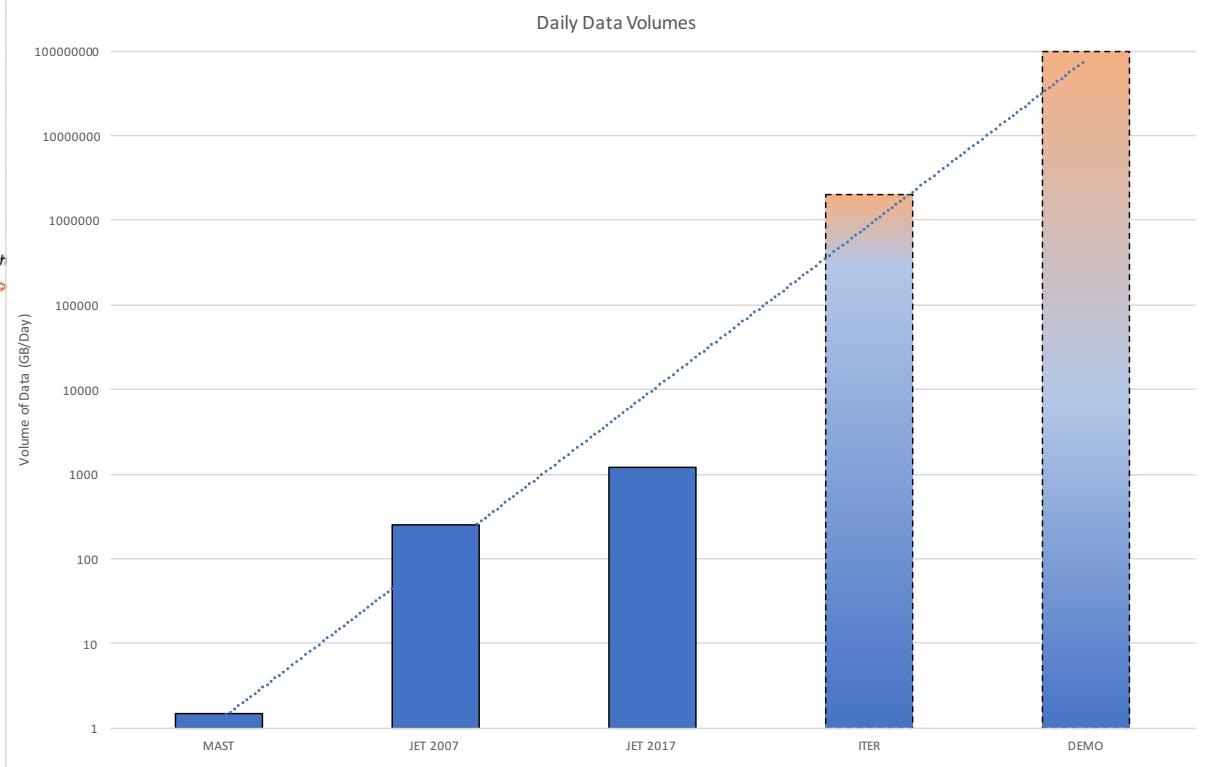
MAST Upgrade – UK industry benefits



Data Challenge

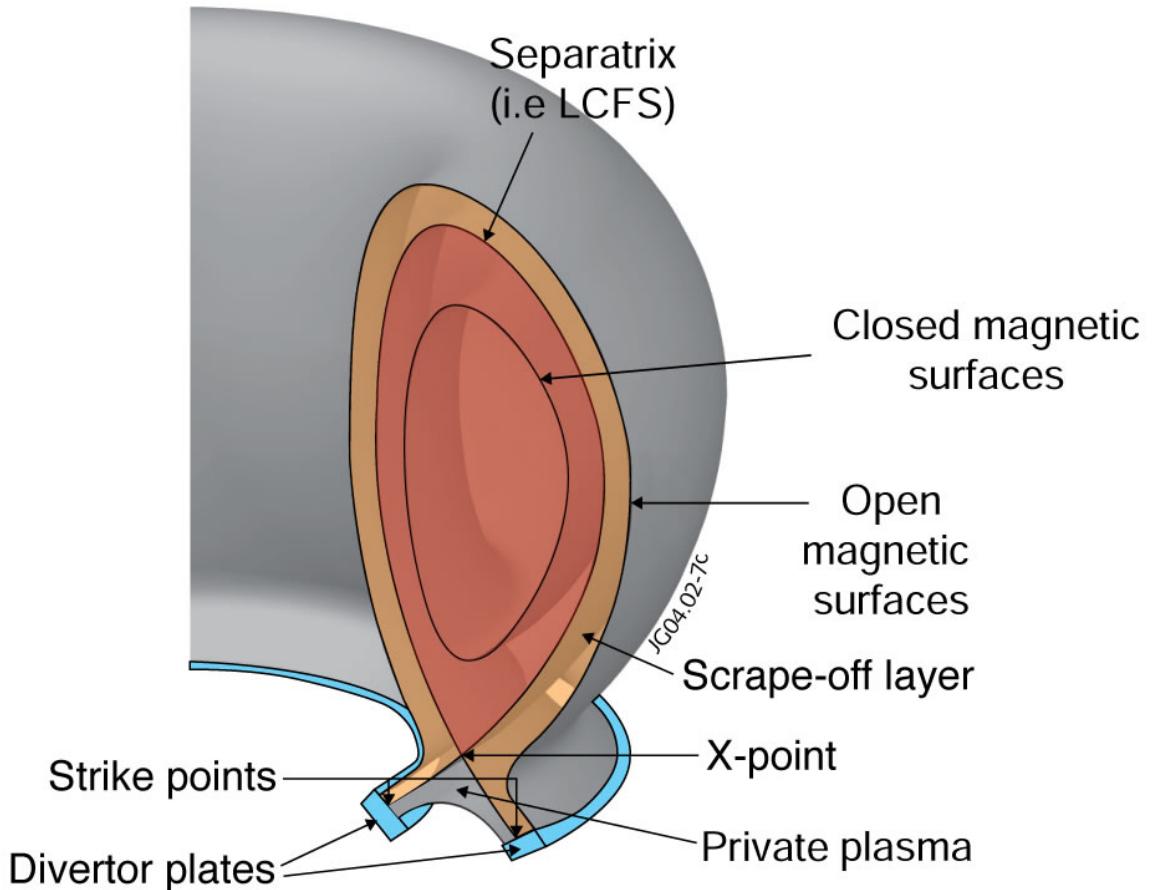


- **ITER is Fusion's First Global Experiment**
 - EU, Japan, USA, India, Russia, China, Korea
 - Will produce more data in 1 day than the total generated by JET in 35 years



Thanks: Shaun de Witt

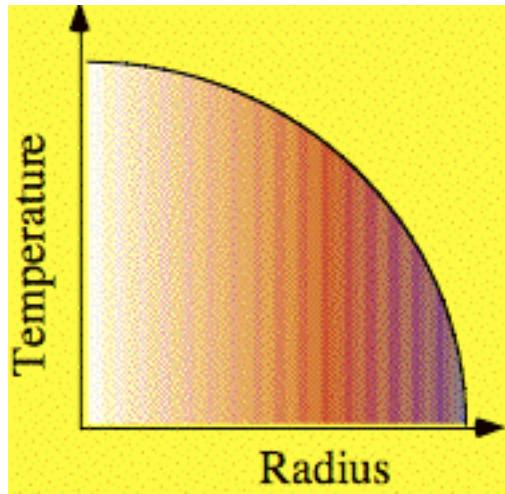
Fusion simulations: multiple



Some examples:

- a) Turbulence
- b) Impurity transport
- c) Instabilities
- d) Plasma exhaust
- e) Energetic particle transport

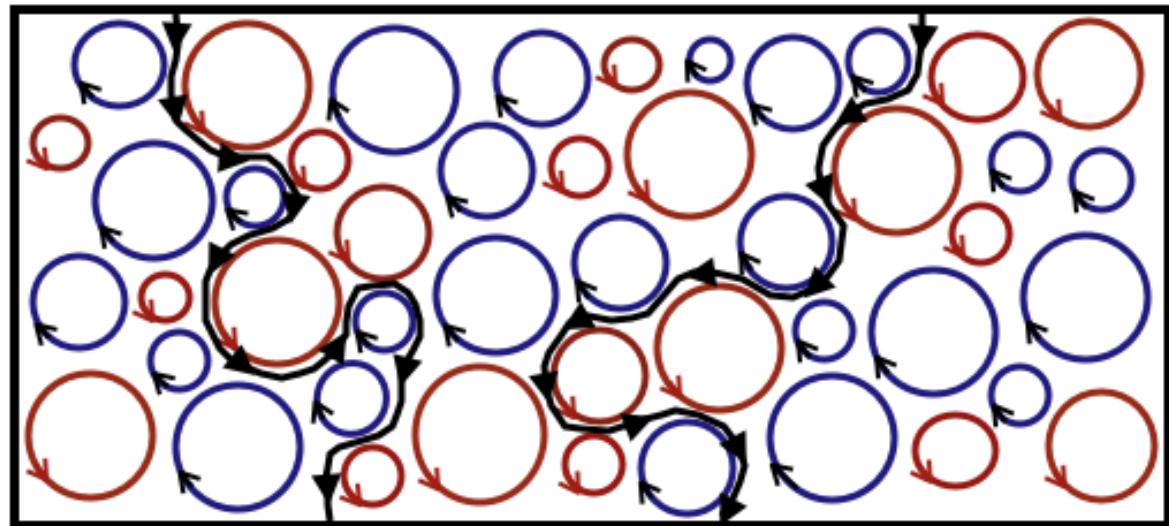
Example problem



Turbulence relaxes the gradient.

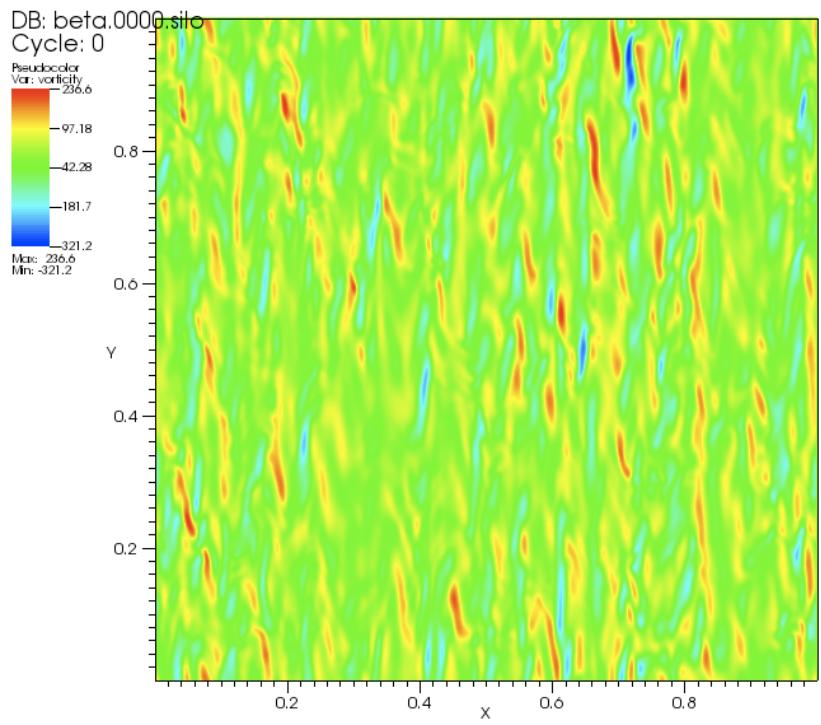
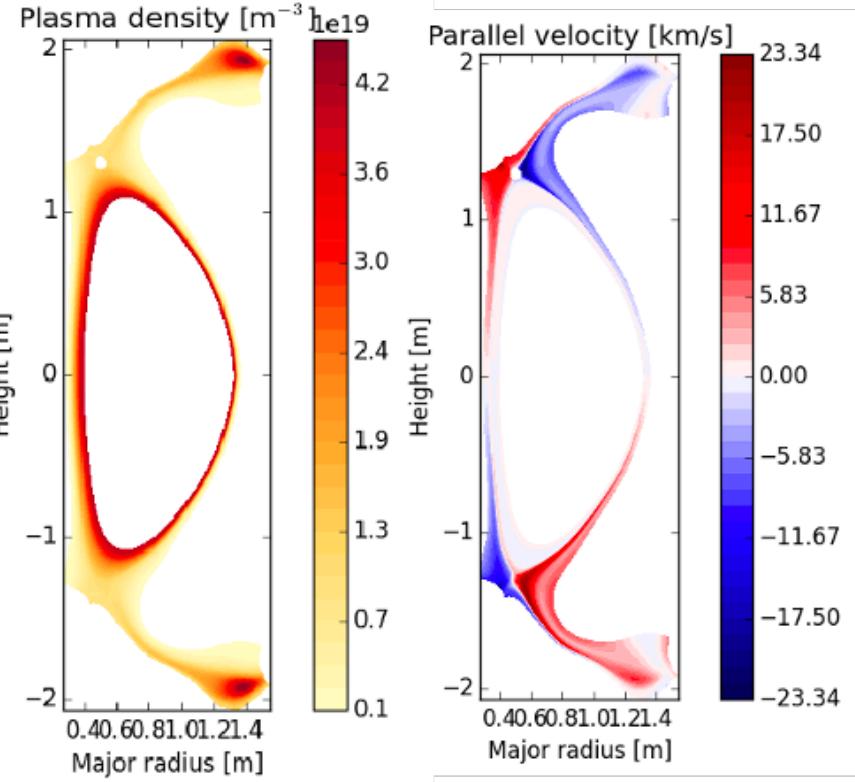
Turbulent transport reduces confinement time.

Hot core, cool walls - but nature abhors gradients



Challenge : TO CONTROL THE TURBULENT TRANSPORT

Turbulent transport in tokamak plasmas is computationally intensive.

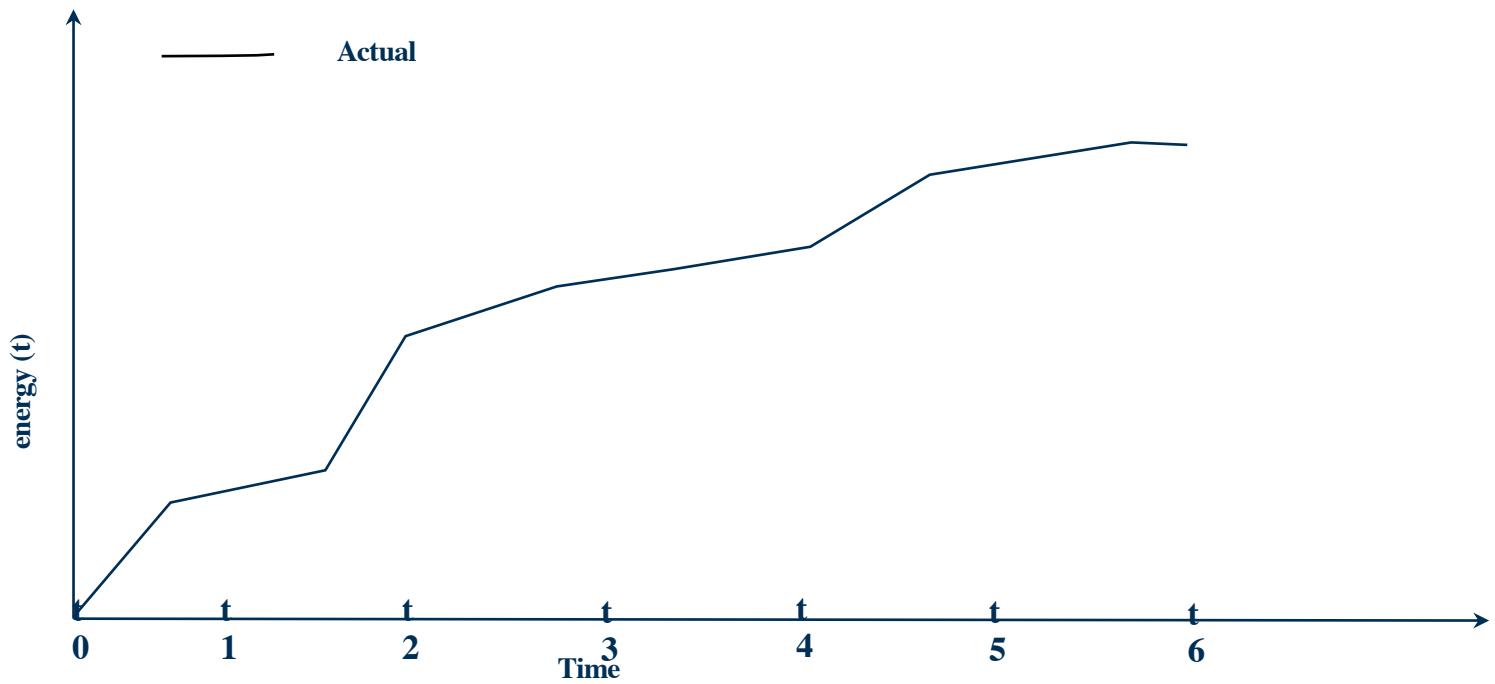


Parallelization options: time

Problem with simulations – extremely long computational time.

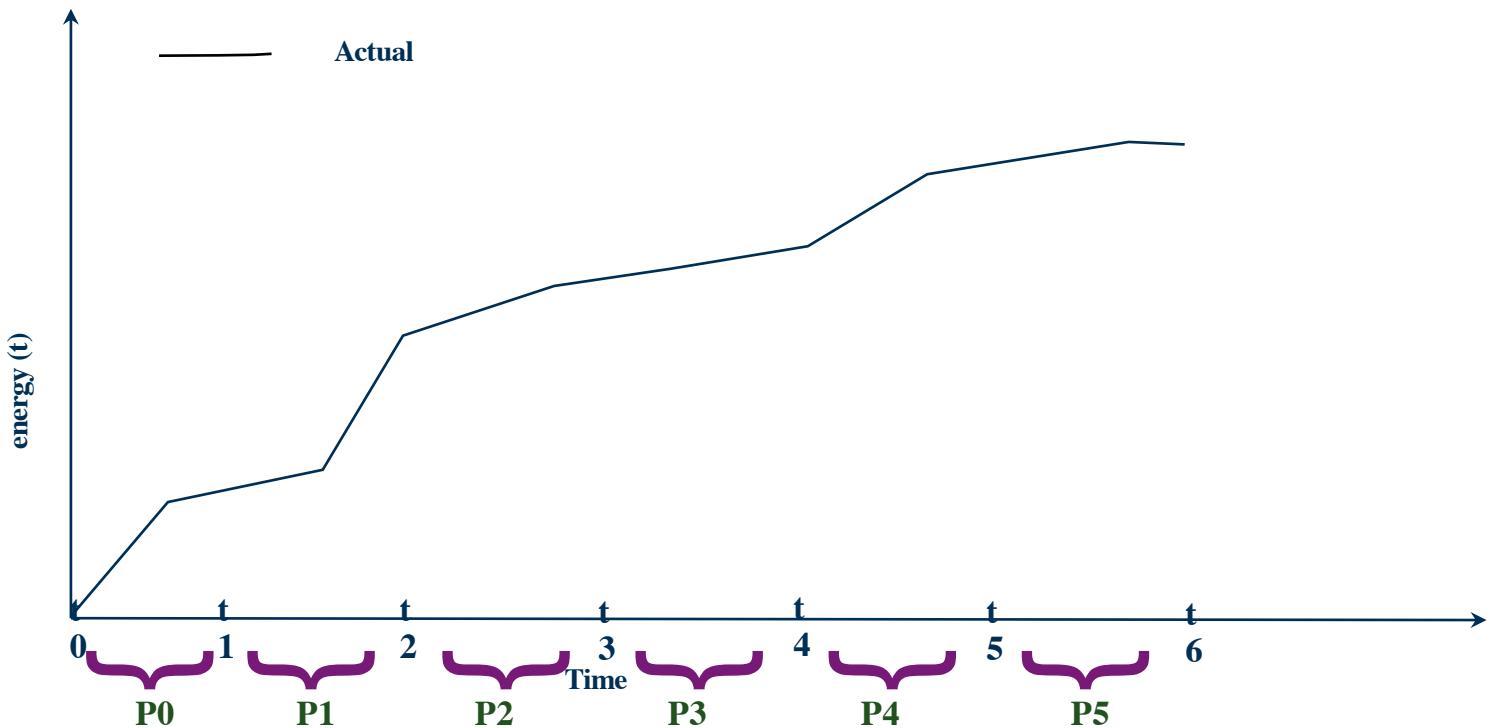
But first let's focus on the solution to achieve the parallel simulation.

F is a propagator evolving the function (energy(t)) from initial time, t_0 , to a later time ...



Parallelize time – solve multiple time slices simultaneously

F is a propagator evolving the function (energy(t)) from initial time, t_0 , to a later time ...



Various algorithms exist for this purpose

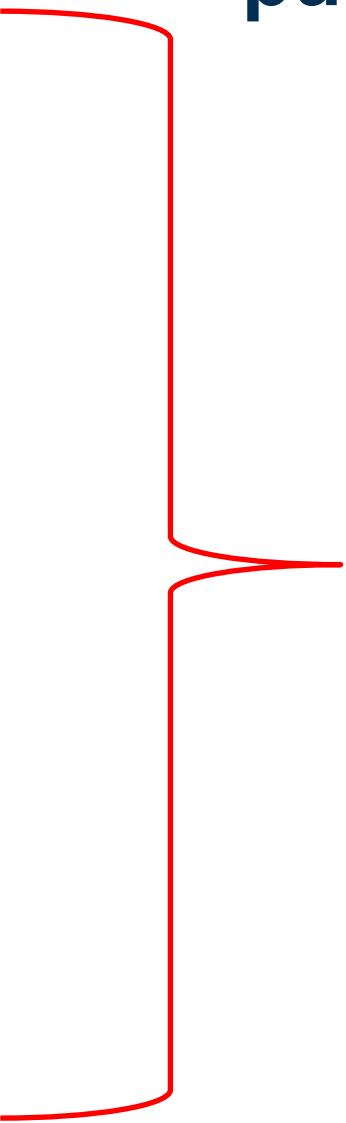
- RIDC
- Parareal
- PFASST

Various algorithms exist for this purpose

- RIDC

- Parareal

- PFASST



Use a predictor-corrector technique

Various algorithms exist for this purpose

- RIDC

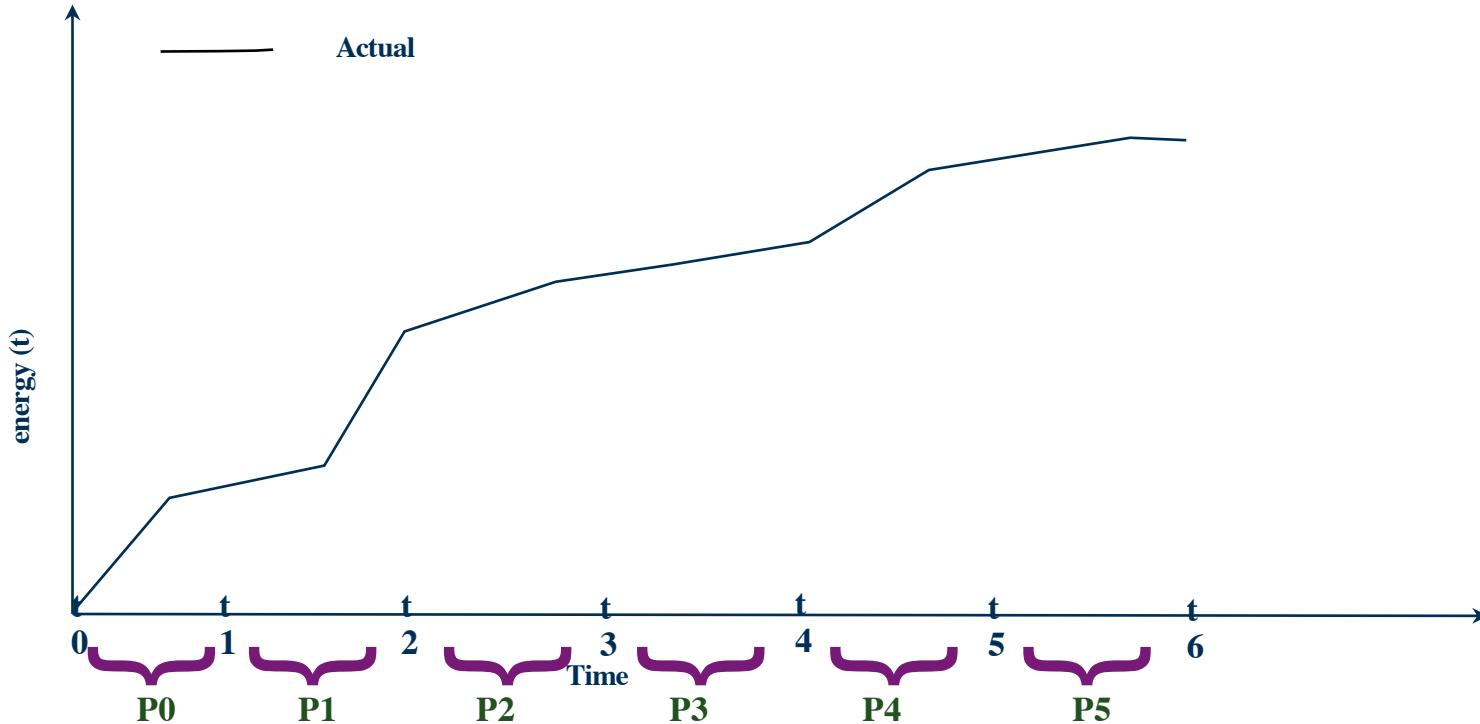
- Parareal

- PFASST

Use a predictor-corrector technique

Overview of the Parareal Algorithm ...

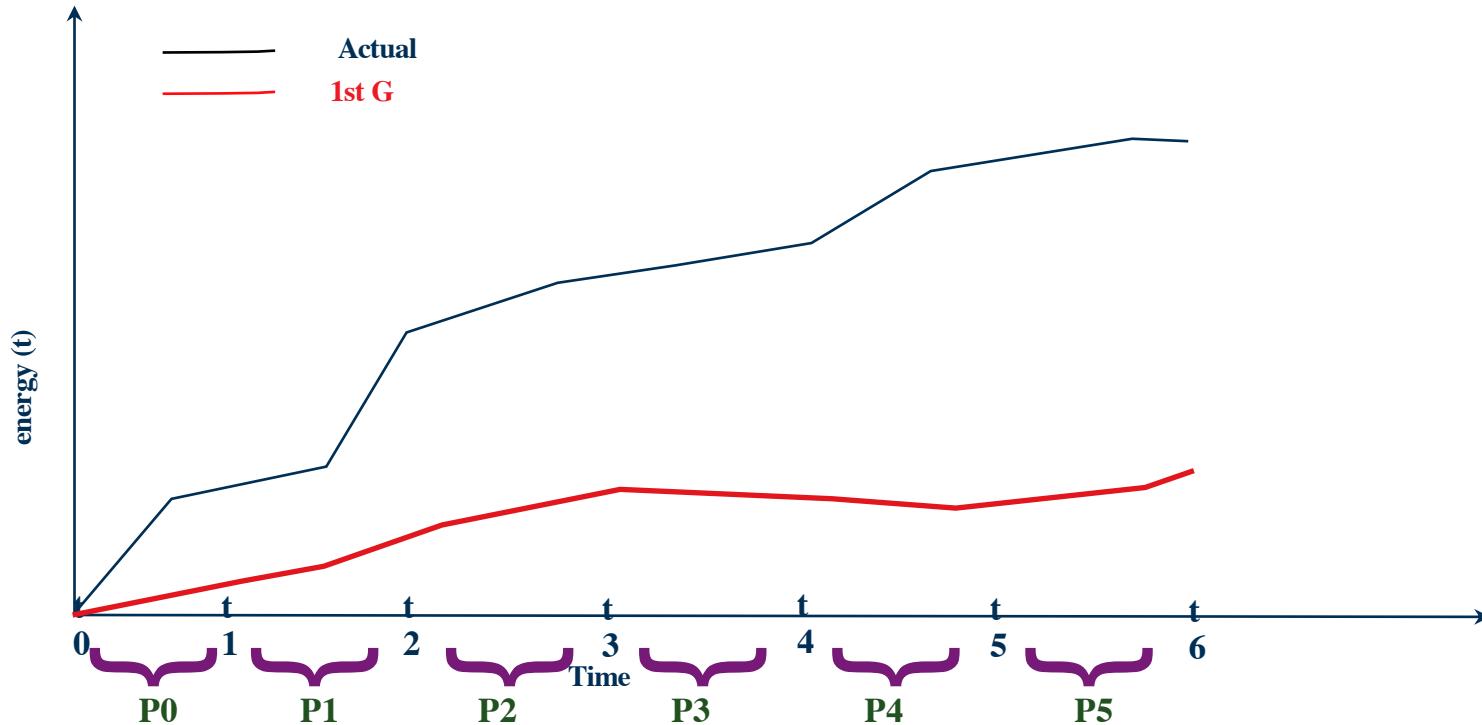
F is a propagator evolving the state, λ . The function, energy(λ, t) thus changes from initial time, t_0 , to a later time ...



$$\lambda_{i+1}^{k+1} = G_{\Delta t}(\lambda_{(i)}^{k+1}) + F_{\Delta t}(\lambda_{(i)}^k) - G_{\Delta t}(\lambda_{(i)}^k)$$

F is a propagator evolving the state, λ . The function energy(λ, t) thus changes from initial time, t_0 , to a later time ...

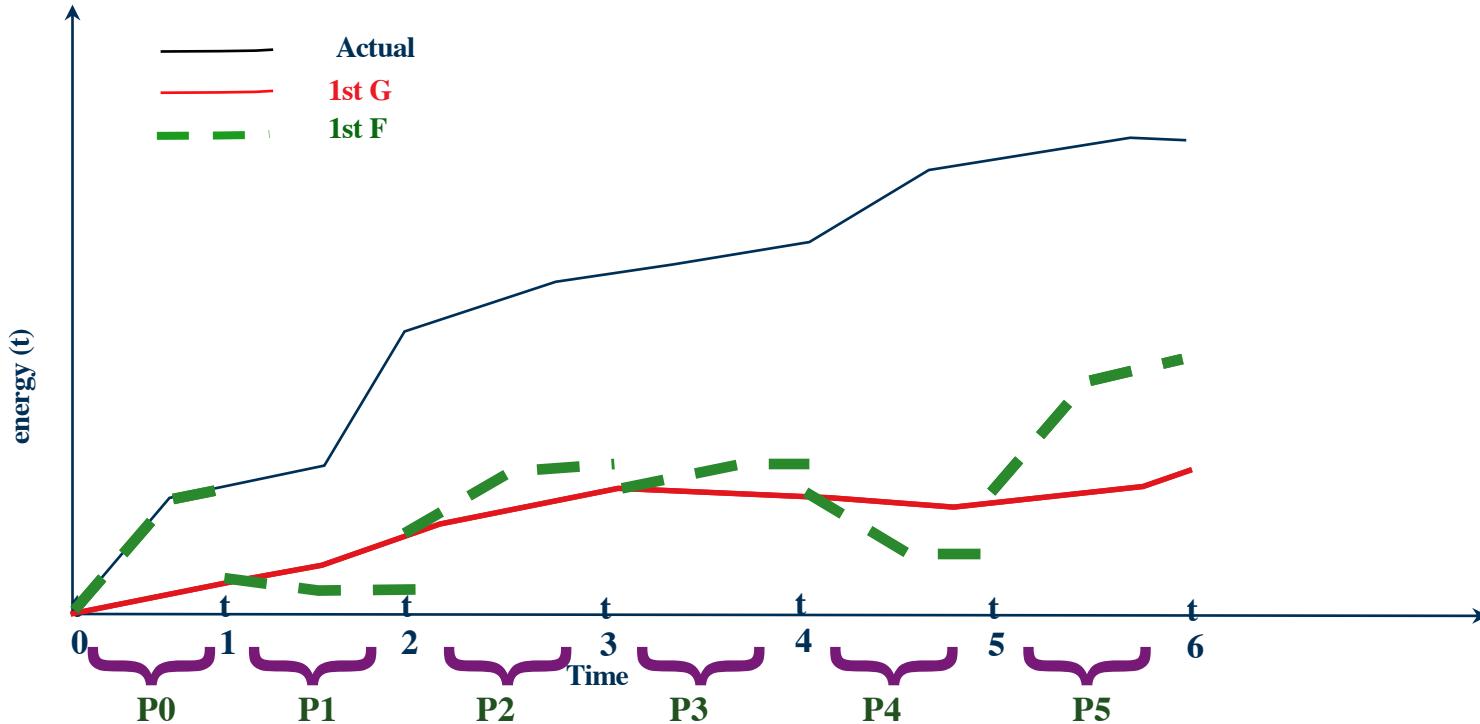
G - faster but inaccurate propagator



$$\lambda_{i+1}^{k+1} = G_{\Delta t}(\lambda_{(i)}^{k+1}) + F_{\Delta t}(\lambda_{(i)}^k) - G_{\Delta t}(\lambda_{(i)}^k)$$

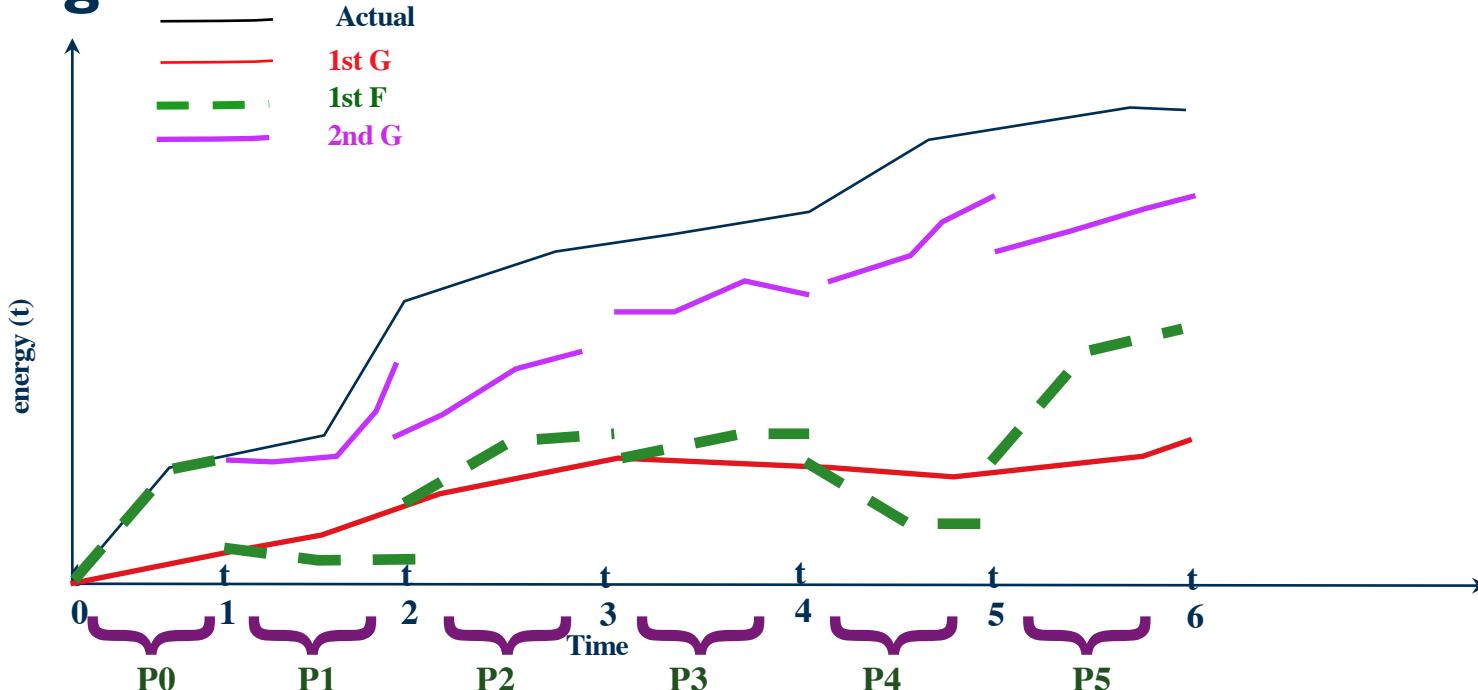
F is a propagator evolving the state, λ . The function, energy(λ, t) thus changes from initial time, t_0 , to a later time ...

G - faster but inaccurate propagator



$$\lambda_{i+1}^{k+1} = G_{\Delta t}(\lambda_{(i)}^{k+1}) + F_{\Delta t}(\lambda_{(i)}^k) - G_{\Delta t}(\lambda_{(i)}^k)$$

F is a propagator evolving the state, λ . The function, $\text{energy}(\lambda, t)$ thus changes from initial time, t_0 , to a later time ... G - faster but inaccurate propagator. Solvers G & F alternate.



$$\lambda_{i+1}^{k+1} = G_{\Delta t}(\lambda_{(i)}^{k+1}) + F_{\Delta t}(\lambda_{(i)}^k) - G_{\Delta t}(\lambda_{(i)}^k)$$

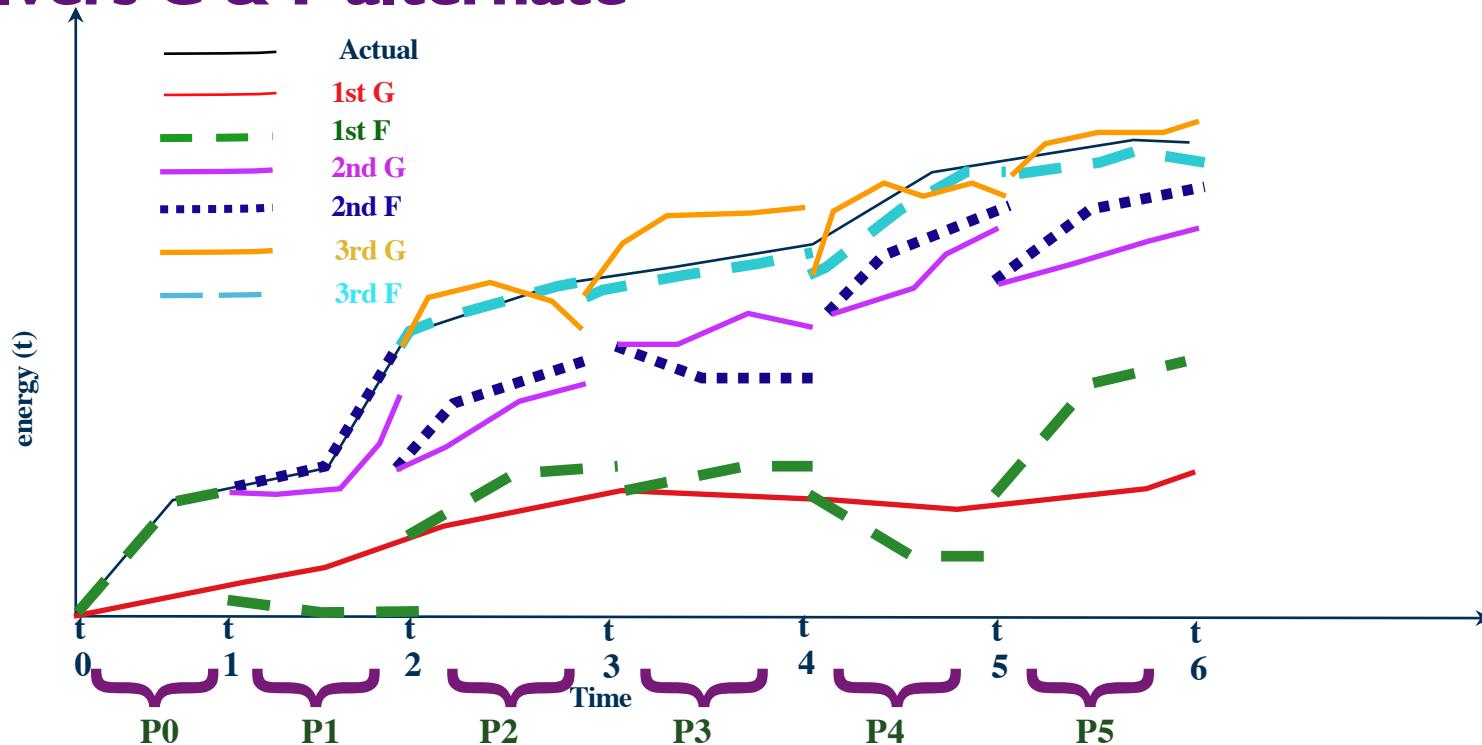
New G

Old G

F is a propagator evolving the function (energy(t)) from initial time, t_0 , to a later time ...

G - faster but inaccurate propagator

Solvers G & F alternate



$$\lambda_{i+1}^{k+1} = G_{\Delta t}(\lambda_{(i)}^{k+1}) + F_{\Delta t}(\lambda_{(i)}^k) - G_{\Delta t}(\lambda_{(i)}^k)$$

↓ ↓

New G **Old G**

Success of Algorithm Depends on Multiple Factors

Algorithm always converges if $k=N$.

★ But, success in achieving significant speedup if $k \ll N$

★ G is much cheaper than F.

★ “Good” G: Solutions converge

★ “Bad” G: Solutions diverge

★ No “fixed recipe” for G !

✓ Despite solutions being very sensitive to initial conditions for - it is possible to choose G.

Selecting Optimum Coarse Solver is Important

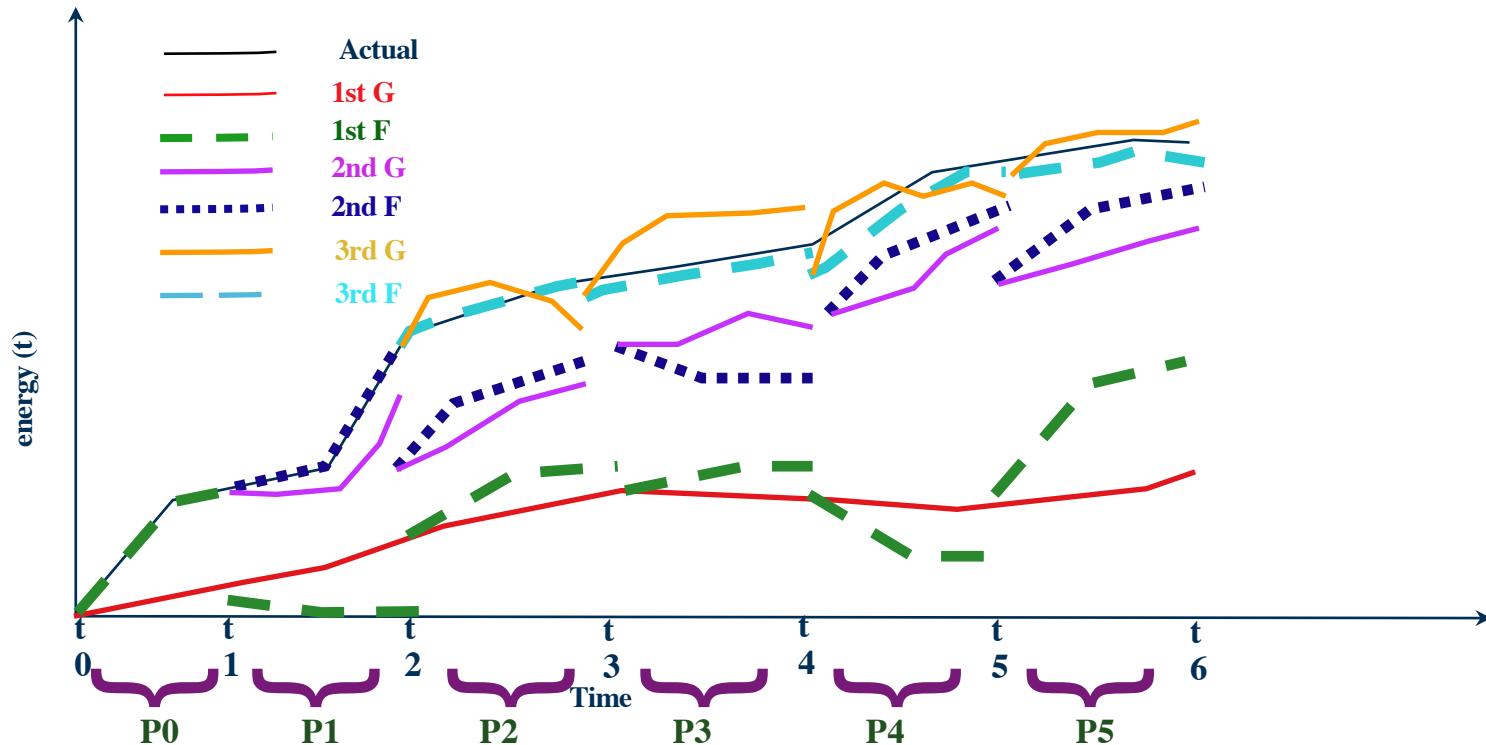
Different approaches can/should be explored to find G. Any one of them, or a combination of them, may work :

- Some of the physics may be ignored when solving with G, to achieve speedup.
- G can be same as F, but may be solved over a coarser k-mesh (or spatial grid).
- G may be same as F, but may be solved with a larger timestep (dt) and less accuracy.
- Use a different G.

Advantages of using the IPS Framework

- **portable parareal framework (L.Berry, W. Elwasif, ORNL)**
- **written in python.**
- **exploring multiple cases with relative ease.**
- **hybrid parallelization (space + time).**
- **Less focus on numerics of parareal scheme.**
- **Prime focus on coarse solver.**
- **Reuse of processors already having attained convergence.**

Event driven parareal dramatically improves computational gain



$$\lambda_{i+1}^{k+1} = G_{\Delta t}(\lambda_{(i)}^{k+1}) + F_{\Delta t}(\lambda_{(i)}^k) - G_{\Delta t}(\lambda_{(i)}^k)$$

↓
New G ↓
Old G

But, what about managing the data?

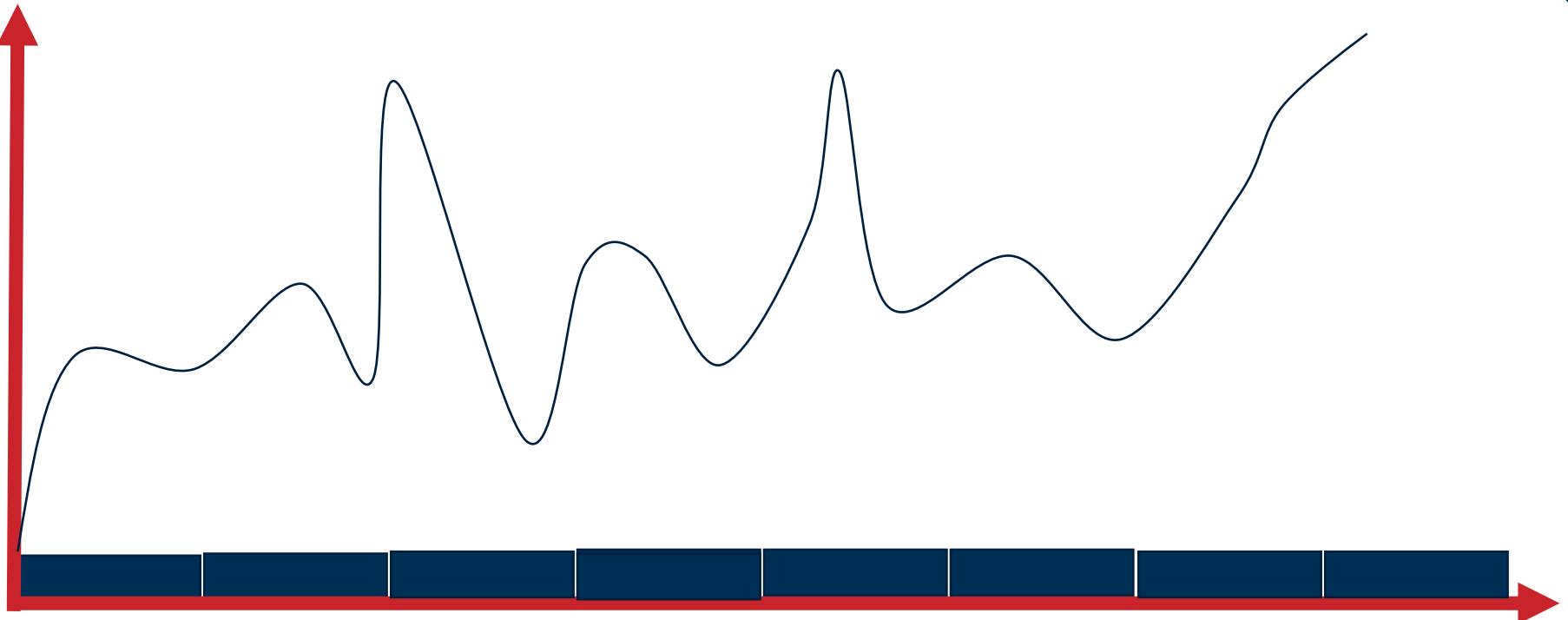
Data storage should provide ease in:

- Visualization
- Checkpointing
- Analysis of random variables

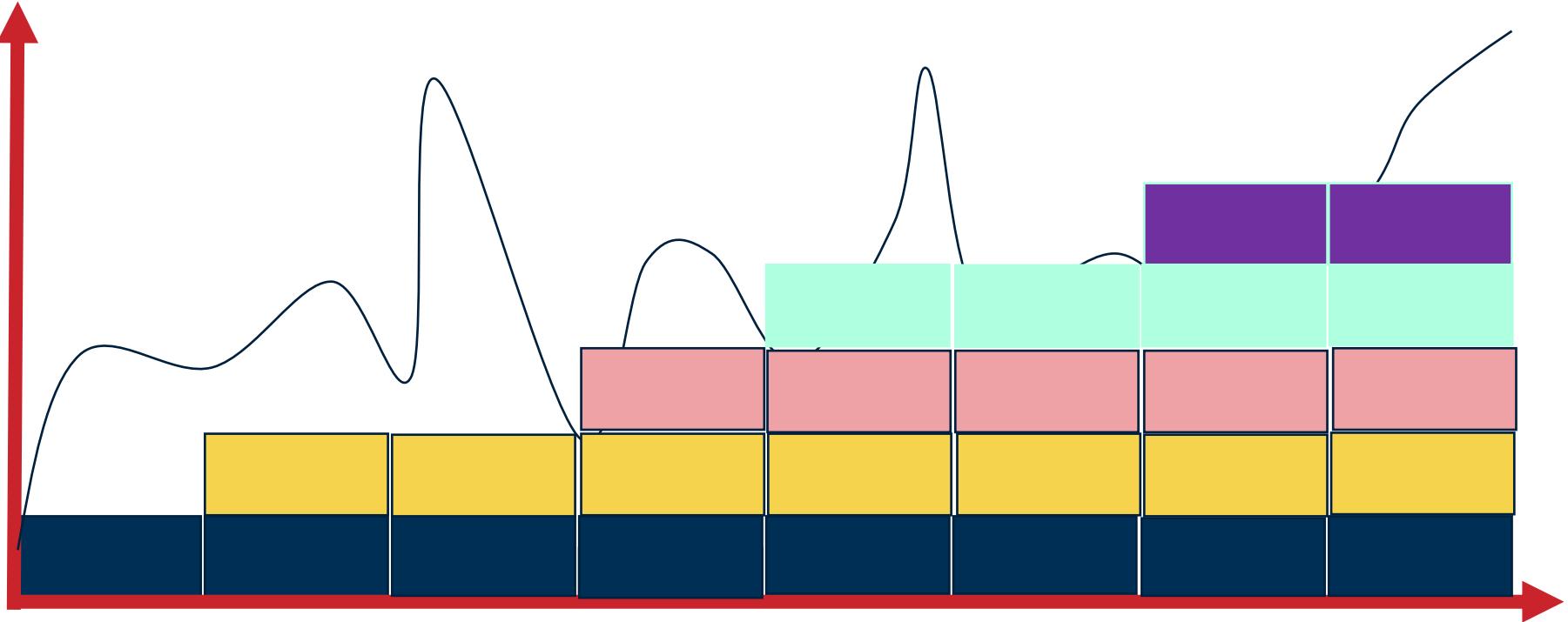
Data size is relatively manageable on current small runs



Data size starts to become difficult with longer runs



Data size is a challenge with parareal runs



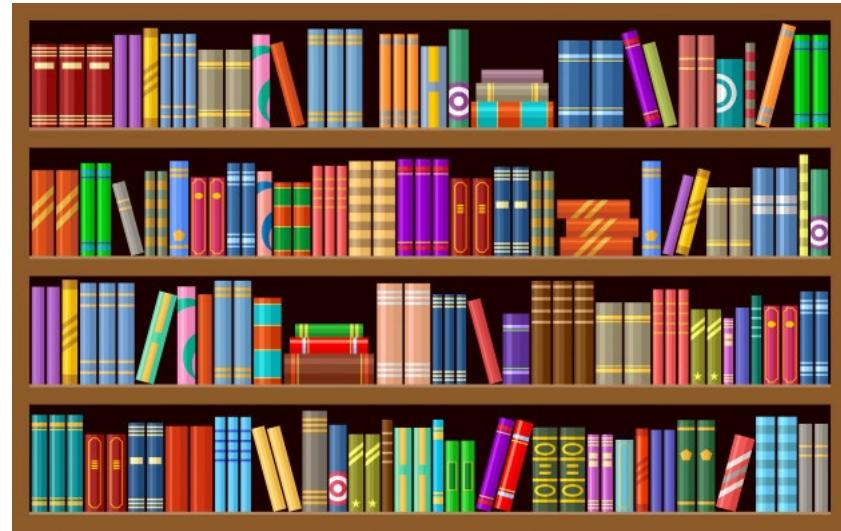
Data challenge in the Exascale

A Parareal run will generate approximately 3 times data per timestep (G+F+corrected value).

Currently, about 100MB per timestep on a serial run. So, 100K timesteps will need 30TB using parareal (per iteration).

For complete physics- 3TB per timestep (So, 300000 TB per iteration). So, when talking of Parareal iteration, k=20(say), we are discussing Exabites.

Current serial run: data retrieval is like looking for a sentence in a given page in a book in a book case.



Exascale simulation: Find the same page with the sentence hidden in a hole somewhere in a BIG city



You first need a map!

How fast can you traverse through the map?!

Code snippet for CORTX I/O

```
file_in=$1
file_out=$2
file_in_data=$1"_FineRun_data"
file_out_data=$2"_data"

id=$(awk 'NR==1{print $1}' $file_in )
echo "id for read in=" $id
addrl=$(awk 'NR==2{print $1}' $file_in)
addrh=$(awk 'NR==2{print $2}' $file_in)

/home/users/jusers/samaddar1/sage/Clovis/clovis-sample-apps_Sept2020/clovis-
sample-apps/c0ct $addrl $addrh $file_in_data 4 $id

./2d_Cond_F $file_in_data $file_out_data
```

... Code snippet for CORTX I/O

```
echo "Done F. Now c0cp output file and save ID & addresses"
```

```
c=$RANDOM  
d=$RANDOM
```

```
/home/users/jusers/samaddar1/sage/Clovis/clovis-sample-apps_Sept2020/clovis-  
sample-apps/c0cp $c $d $file_out_data 4 -a 8 | awk '{print $2}' > $file_out
```

```
echo $c $d >> $file_out
```

Steps for porting code to tiered storage system.

1.

Legacy serial
code

2.

Code in
Parareal
framework

3.

Parareal version on
Cortex

References:

- Lions J. L et al, C.R. Acad. Sci. Paris, Serie I, 332 (2001), pp.661668
- Samaddar D. et al, J. Comp Phys, 18 (229) (2010)
- Berry, L. A. et. al, Journal of Computational Physics 231(2012) 59455954
- Elwasif W.R et al, 4th IEEE Workshop on Many-Task Computing on Grids and Supercomputers, MTAGS 2011, (2011)
- Samaddar D. et al, Comp. Phys. Comm. 221(2017), pp.19-27

Thank you!