

# I/O IMPROVEMENTS AT NCEP

George VandenBerghe  
Lynker Technologies  
NOAA/NWS/NCEP/EMC

October 13, 2023

# Overview

- NCEP is the United States' premier NWP institution.
- First formed in 1954.
- Available compute of that time was a few kiloflops.
- Today it exceeds 700 teraflops. (actually achieved by applications and small fraction of 13PF theoretical peak)

# NCEP COMPUTING

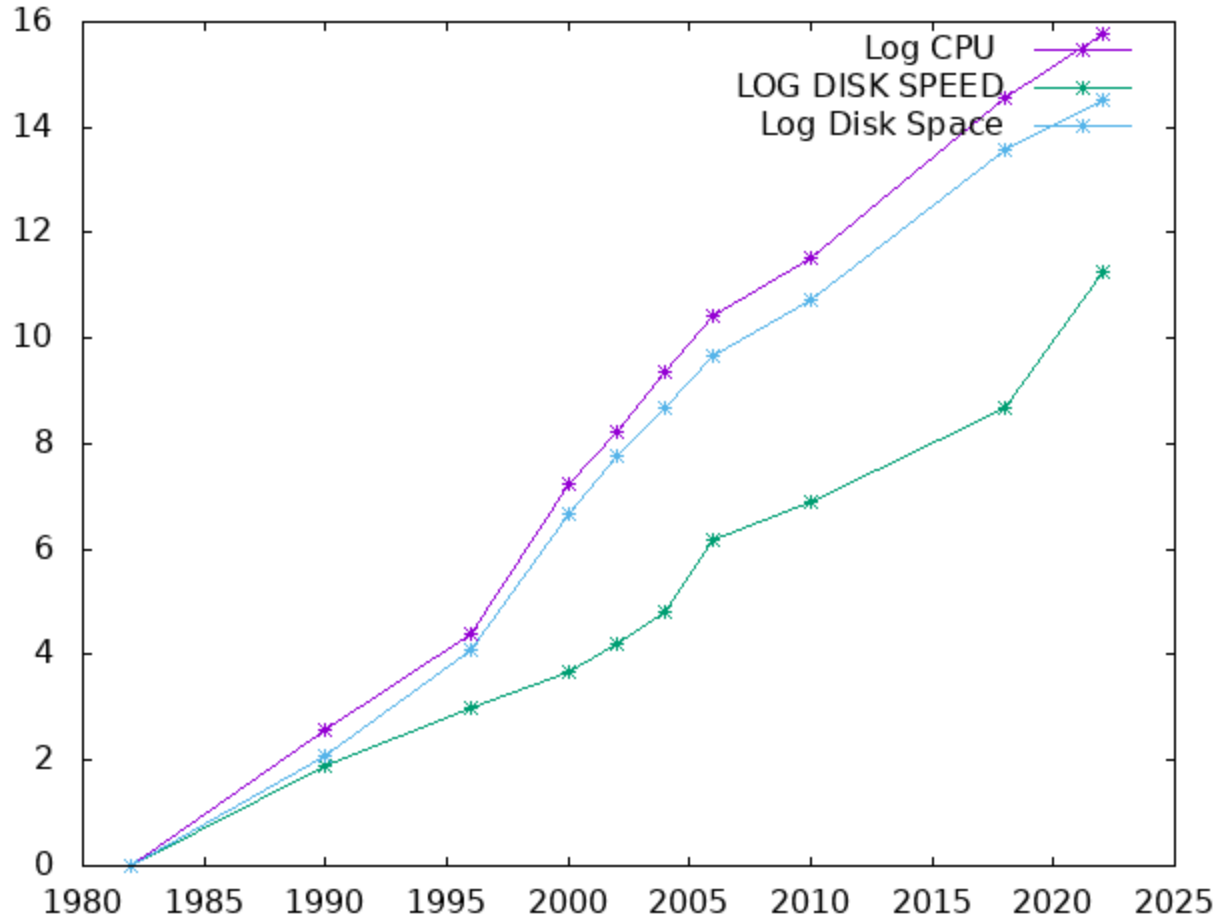
- There have been ~~28~~ ~~31~~ ~~32~~  
**38** doublings in compute capacity since 1954.
- The crossed out figures were from presentations in 2000 and 2006-8
- .

# Past Platforms.

- ENIAC 1 kflop
- IBM 70x mid 50s 10 kflop
- IBM709x early 60s 100 kflop
- CDC6600 mid 60s-early 70s: 1 mflop
- IBM 360/195 x3 70s-early 80s 10 Mflop
- CDC CYBER205 x2 80s 100 Mflop
- CRAY Y-MP8 early 90 1Gflop (1.3 aggregate)
- CRAY C90 mid to late 90s 6 Gflop (8 aggregate)
- IDM SP 1999-2000 30 gflop (70 aggregate)
- IBM SP 2000-2002 60 gflop (140 aggregate)(x2)
- P690 (P4) 2003-2004 160 gflop (370 aggregate)(x2)
- P655 (P4+) 2005-2006 500 gflop (1150 aggregate)(x2)
- P575 (P5) 2007-2008 1500 gflop ( 3400 aggregate)(x2)
- **WCROSS2 2022-??** **700000 gflop**  
(x2)

- ***Available cycles have doubled every two years since 1950 with more rapid doubling in recent two decades*** (correlation with my career is coincidental).

# Compute v.s. disk space and speed



# A FEW SITE TRENDS (2002)

- Compute requirements grow exponentially.
- E-fold time ~2 years.
- Compute budget near constant. (e-fold time 20+ years)
- Watch out for per flop costs that are constant or decreasing slowly ( with e-fold time >2 years)
- (floor space/tflop, disk/tflop, tapes/tflop)
- Disk metrics/compute was exponentially declining to about 2018, but this trend has improved (good!)  
Exponential decline to 2018, but WCOS2 provided a jump.

# Disk Farm Service Metrics.

- 1982, Cyber 205 15 mbytes/second 10GB
- 1990 Cray YMP 100 mbytes/second 80GB
- 1996 Cray C90 300 mbytes/second 600GB
- 2000 IBM/SP 600 mbytes/second 16TB
- 2002 P690 1000 mbytes/second 24TB
- 2004 P655 2000 mbytes/second 60TB
- 2006 P575 8000 mbytes/second 160TB
- WCROSS Intel.(\$xx) 16+32+24, 72Gbytes/second  
13PB (scattered across three clusters)
- WCROSS2 (amd) 1TB/sec 20PB

# I/O Issues (from 2008.. still valid)

- *Disk space/flop has scaled linearly.*
- *Disk performance per flop has not!* (But this has gotten better since 2022 WCOS2!)
- *This long going IT trend was masked by low I/O requirements of traditional NWP.*

*I/O desires are now a linear function of compute capacity.  
( $I \sim k \cdot C$  where  $C$  is compute) (used to be less because of CFL.. no longer true)*

*I/O constraints became noticable in 2006 and painful by 2014.*



# I/O

- Disk metric trends were MUCH flatter than compute. (Both transactions/sec and bandwidth). WCROSS2 provided notable improvement
- Flash metadata revolutionized many tiny file operations (such as modern software builds ). We started seeing this 2015 or so
- Burst buffers, after initial skepticism have also increased IOPS although this is harder to quantify. WCROSS2 goes further with full flash filesystem available. (surprisingly this does NOT improve many apps' I/O )

# And Speaking of Memory

- 1982 C205 32mbytes
- 2008 P575 node 32 gbytes
- 1000x increase in memory
- 100x increase in cpu/node
- Flush memory to single disk 1982 10 sec  
2007 500 sec.
- Flush ½ memories (16G\*160 nodes) to filesystem 1982  
2 sec 2007 320 seconds. 2022(500 seconds)
- (This is a virtuous surprise.. BW is almost keeping up  
with CPU WCROSS2 in particular JUMPED)

# I/O Mitigations

A. Novel techniques, non-filesystem object stores, NVM, on chip caching with separate object handlers .....

1a. NOT DISCUSSED IN THIS TALK (we're looking at it)

**B. JUST .. DO ... LESS .... I/O !!**

That's the topic of the rest of this talk

# Workflow Mitigations

- Typical NCEP pattern was
- Forecast Write → post process → write → product generator → write.
- This generates many writes followed by reads followed by more writes.
- Alternative is to make post processor(done) and perhaps product generator a part of the forecast model. NCEP is at an intermediate stage there. **DONE FOR POST PROCESSOR. NEEDED for BC GENERATORS.**
- Disk I/O is replaced by interconnect transfers which are much faster and cheaper.

# UFS(GFS) Deterministic requirements

- Write hourly state and products for first 128 hours
- Write 3 hourly state to 240 hours.
- write 12 hourly state to 384 hours.
- Hourly writes are done every 12-15 seconds. Each write takes 60 seconds (asynchronously)
- State transfer to write group <1 second.

# Inline Post process savings

Write 31 gbytes - - read 31 gbytes - - write  
3.4 gbytes final products

NEW write 3.4 gbytes final products.. 58  
GB I/O saved

**Remember we do one of these every 12-  
15 wall seconds**

We still need the history file though (but  
we've also dealt with that!)

# History file savings

- History files are 32 bit snapshots of model state.
- These compress (in parallel) about 1.5x
- BUT.. we can tolerate precision loss
- We truncate 14 out of 23 mantissa bits. Result (not obvious) results in about 5x compression reduction. (obviously shouldn't restart from this or do high precision calculations but good enough for BC, their primary remaining use)

# Plumbing improvements

- Serial compression was an order of magnitude too slow.
- NetCDF developers provided (on short notice) a parallel compression API for netcdf RECORDS (using HDF5 compression underneath)... spec to working implementation took five months (Nice work!!)



# More plumbing

- Forecast model progresses a model hour every 15 wallclock seconds. Forecasts are written asynchronously.
- Multiple I/O groups in parallel do compression, inline post, and writing, concurrently. (current operational GFS needs eight but will usually keep up with just five.. eight is for filesystem write variation margin) Each group uses 48 MPI ranks at 24 ranks per node.. (this took some tuning)
- Not unique, many models in many institutions evolved independently in this direction.

# Back to history file

- We still need history files for Boundary Conditions for Regional models.
- Inline BC is hampered by the possibility of desired domains unknown at GFS runtime, so getting rid of the history file isn't likely soon.

# Restart file

- We also write restart files for checkpoint and recovery purposes only.
- No downstream product needs them (**but needed for next cycle data assimilation**)
- Restarts are full precision full state. No compression here!
- Work to also parallelize restart writing is ongoing .. several good options available and this is not at all intractable.

# Summary

- After gross declines 1990-2018, I/O metrics are now (almost) keeping up with cpu. But site was I/O bound in 2018 and we're still approximately at that level of constraint in the hardware. **WCOS2, provided big boost.**
- NCEP is doing MUCH less I/O in the GFS, This makes the constraint above.. much less “constraining”
- We do it by reducing precision requirements in history file, thus compressing it 5x, and doing high precision post requirements in the model rather than from state. 18X I/O reduction for primary products.

Questions??