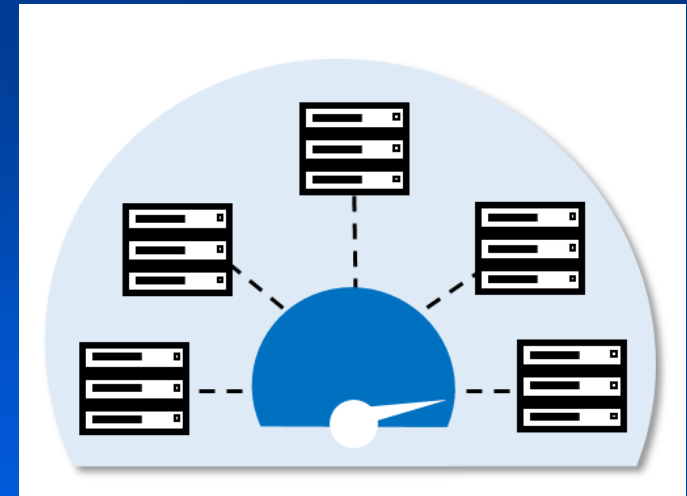


CDRH HPC Clusters



CDRH HPC Team

Stuart Barkley

Dillip Emmanuel

Rusif Eyvazli

Fu-Jyh Luo

Mike Mikailov

Nicholas Petrick

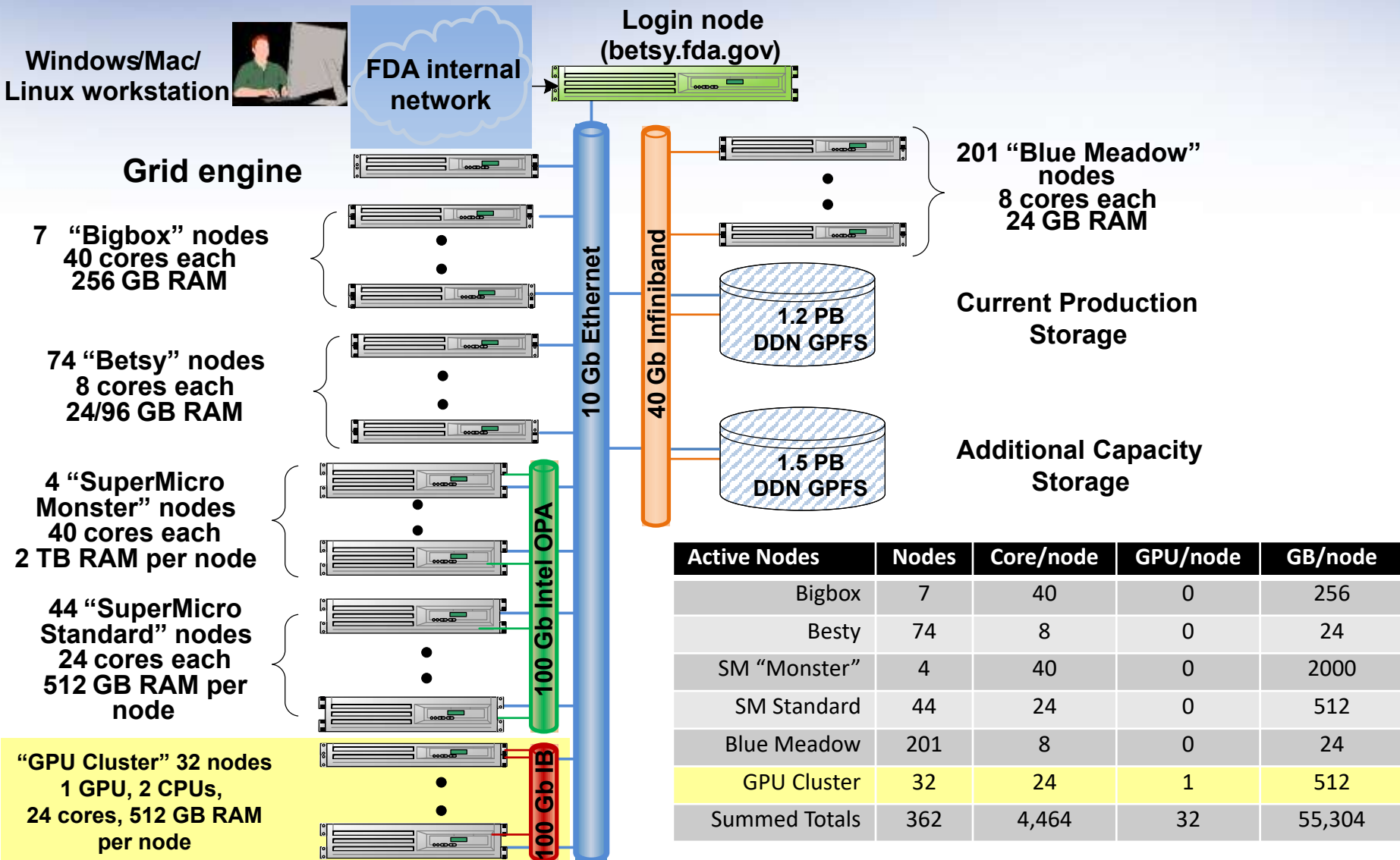


HPC Computational Environments

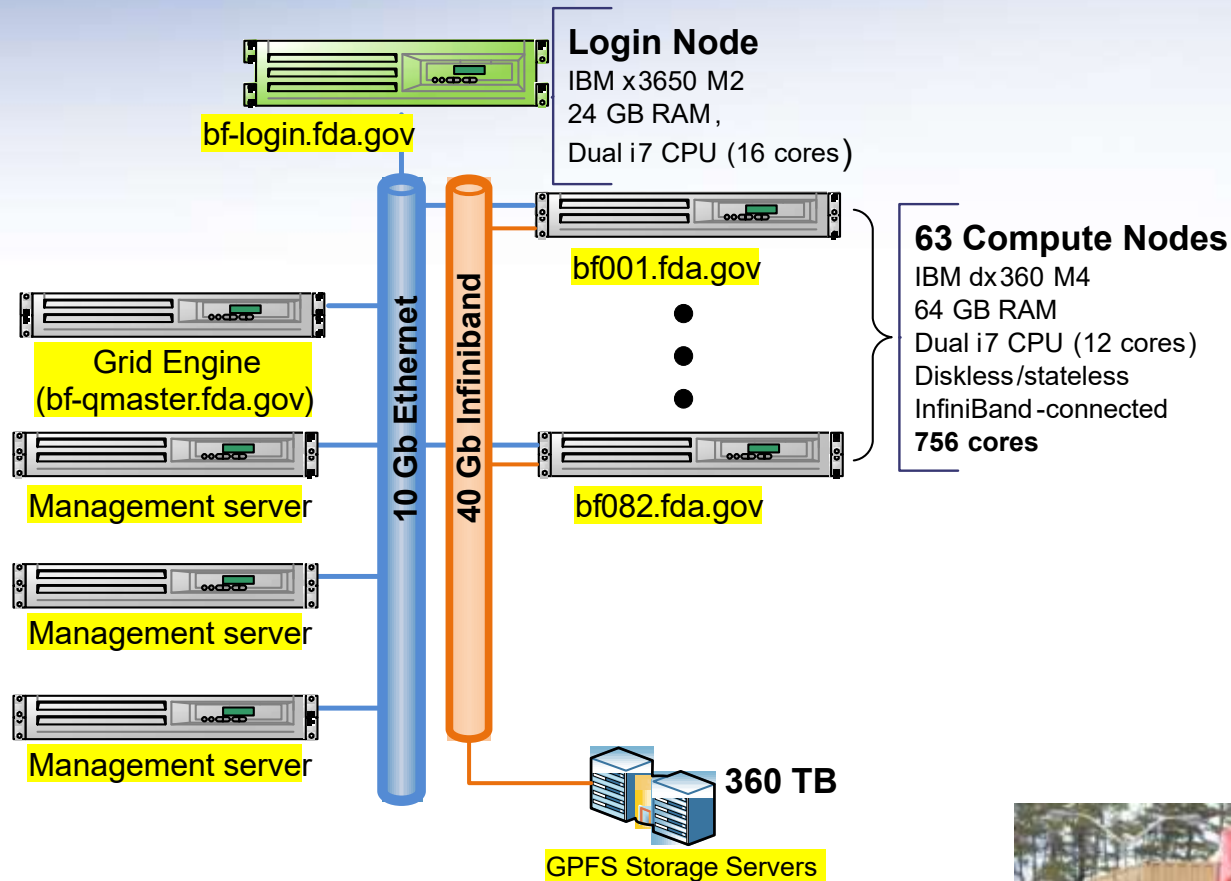
- CDRH Betsy/Bluefin Environments
 - General purpose HPCs to support a wide array of projects
 - Web-based bioinformatics analysis (Galaxy)
 - Artificial intelligence/machine learning
 - Genomics, next-generation sequence analysis and alignment,
 - Modeling and simulation
 - Statistical analysis and more



Betsy Cluster



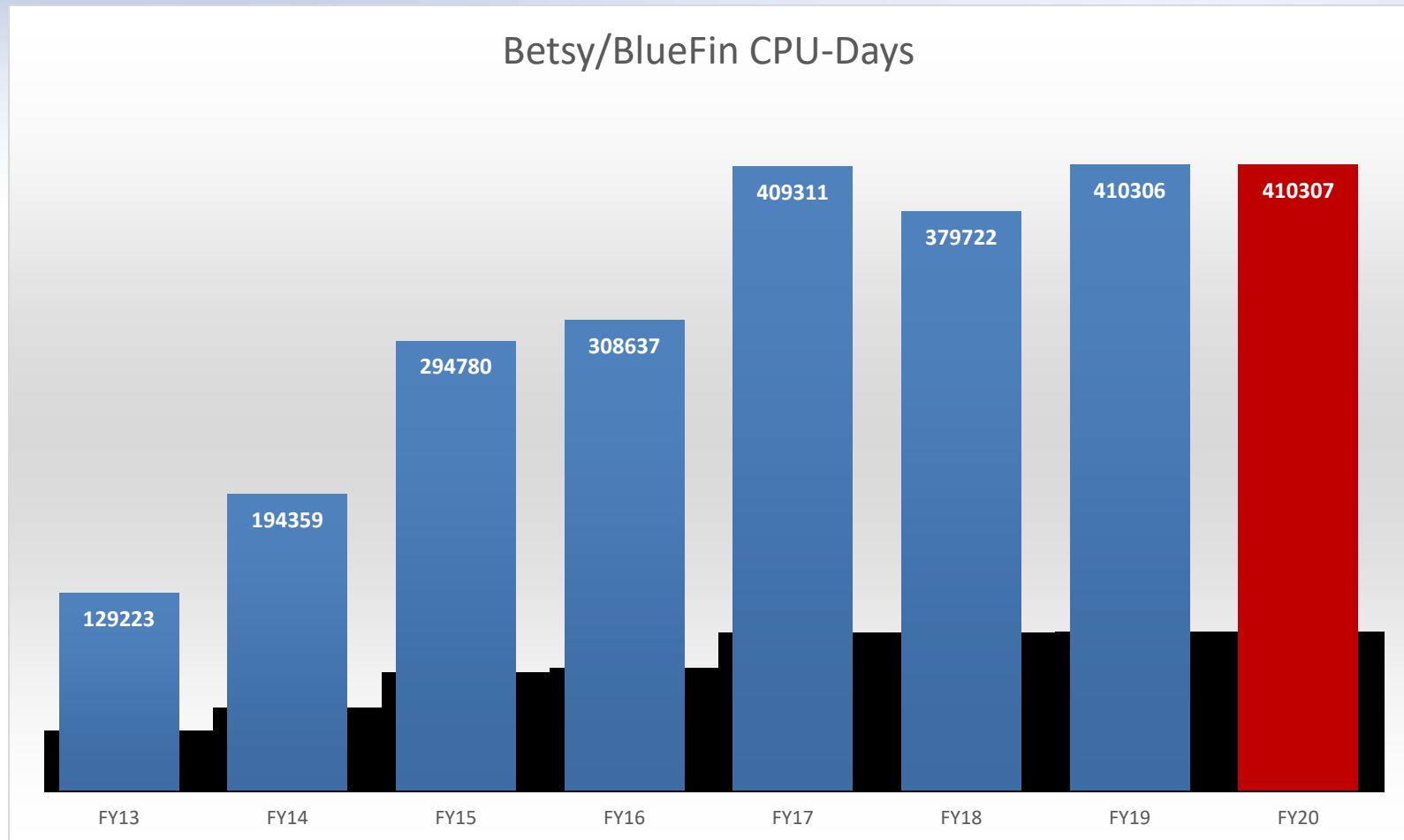
Bluefin Cluster



- Bluefin is a “mobile” computing platform
 - Currently located at CVM facility in Beltsville



Betsy/Bluefin Utilization



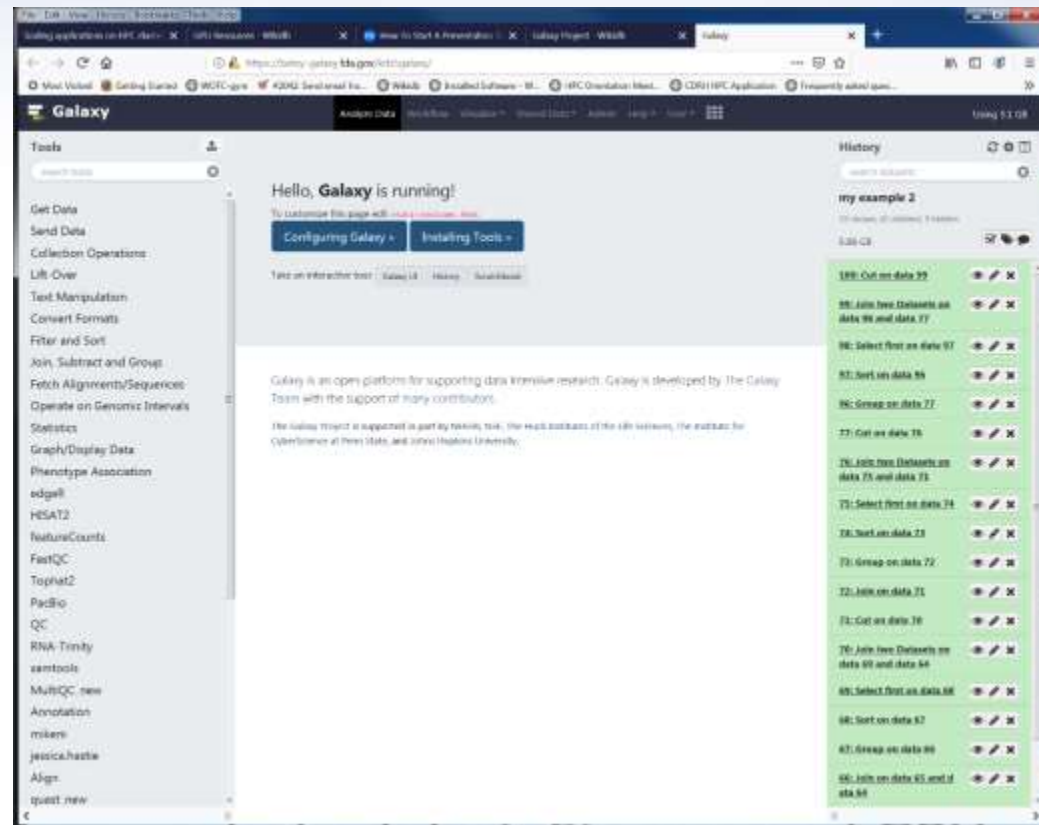
Scientific Applications

300 FDA-approved applications:

- **Computational Chemistry:** Amber, Autodock, Desmond, LAMMPS, ...
- **Computational Fluid Dynamics:** OpenFoam
- **Image Analysis:** CellProfiler, Freesurfer, FSL, ImageJ ...
- **Linkage/Phylogenetics:** CD-HIT, GARLI, Mothur, PHYLOSHOP, POY, QIIME ...
- **Mass Spectrometry:** TPP (PeptideProphet, ProteinProphet, ASAPRatio, ...)
- **Mathematical/Statistics/Modeling & Simulation:** R, SAS, Chaste, FluTE, GridMathematica, ...
- **Next-Generation Sequencing:** ABySS, ALLPATHS-LG, Bcl2fastq, Bedtools, Bowtie, Bowtie2, Bwa, CCMpred, Cufflinks, GATK, GNUMAP, HIVE, Kraken, Mira, Picard, Samtools, Tophat
- **Sequence Analysis & Alignment:** Blast, Blat, Cortex, Dfam, Exonerate, Geneious, GotoBLAS2, HMMER, Jalview, JELLYFISH, Mauve, mpiBLAST, Mugsy, MUMmer, MUSCLE, NovoAlign, ParSNP, Pfam, PRINSEQ, ScalaBLAST, SGA Assembler, SNP Pipeline, SPAdes, SUPERFAMILY, Tablet
- **Multi-paradigm environments:** Parallel MATLAB, Octave, R, Galaxy, ...

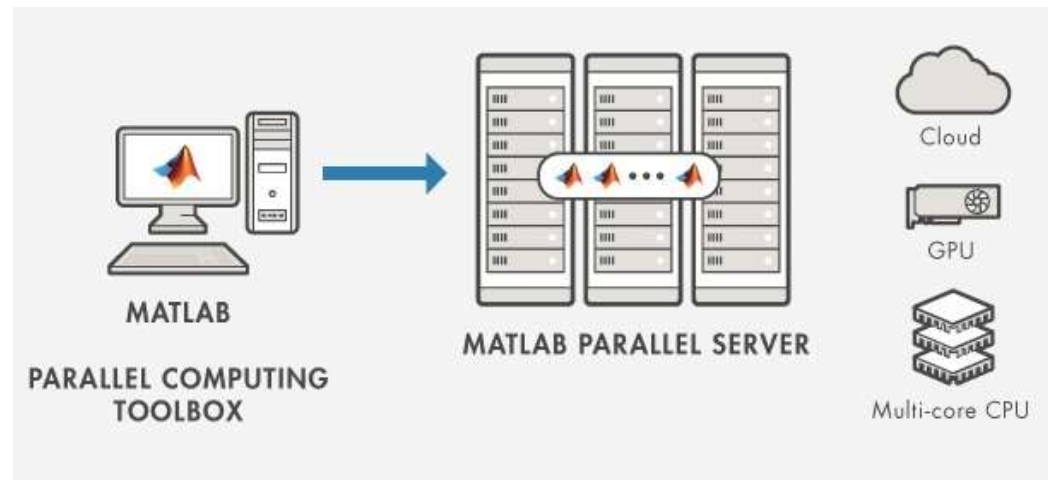
Galaxy Web Annalistic Platform

- Open-source platform
 - Enables researchers without informatics expertise to perform computational analyses through web
 - Users upload data and define/run analysis pipelines
 - Galaxy Tutorials
 - [Galaxy 101](#). How to create workflows.
 - [Uploading data](#) - How to get data into Galaxy.
 - [Learn Galaxy](#)



MATLAB Parallel Server

- MATLAB Parallel Server
 - Lets users scale MATLAB[®] and Simulink[®] programs to HPC cluster
 - Runs programs/simulations as scheduled applications on cluster
 - Desktop license profile dynamically enabled on cluster, so no need to supply MATLAB licenses for cluster



HPC Application Scaling



Mike Mikailov
CDRH/OSEL/DIDSR

Process Scaling

Before parallelization
Single run

```
max=2000
for ( i in 1:max)
{
  [computations]
}

[summarization]

[final results]
```



After parallelization
Many parallel runs

```
i=1
[computations]

[partial results-1]
```



...

```
i=2000
[computations]

[partial results-2000]
```



Concatenation,
summarization
within minutes

```
[partial results
concatenation]

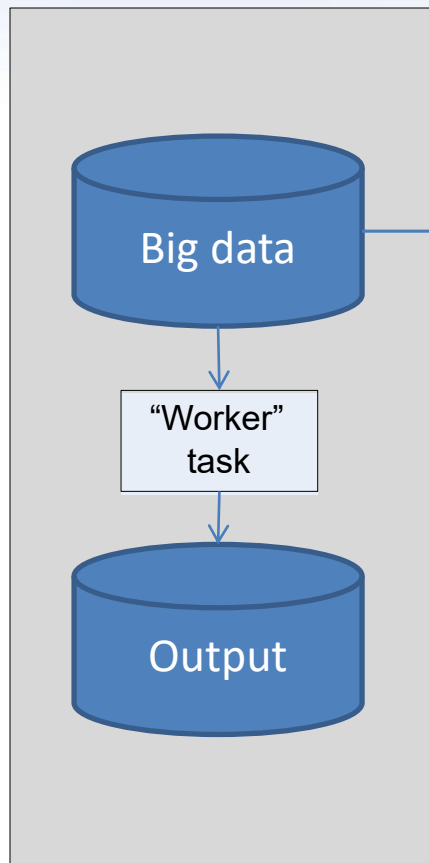
[summarization]

[final results]
```

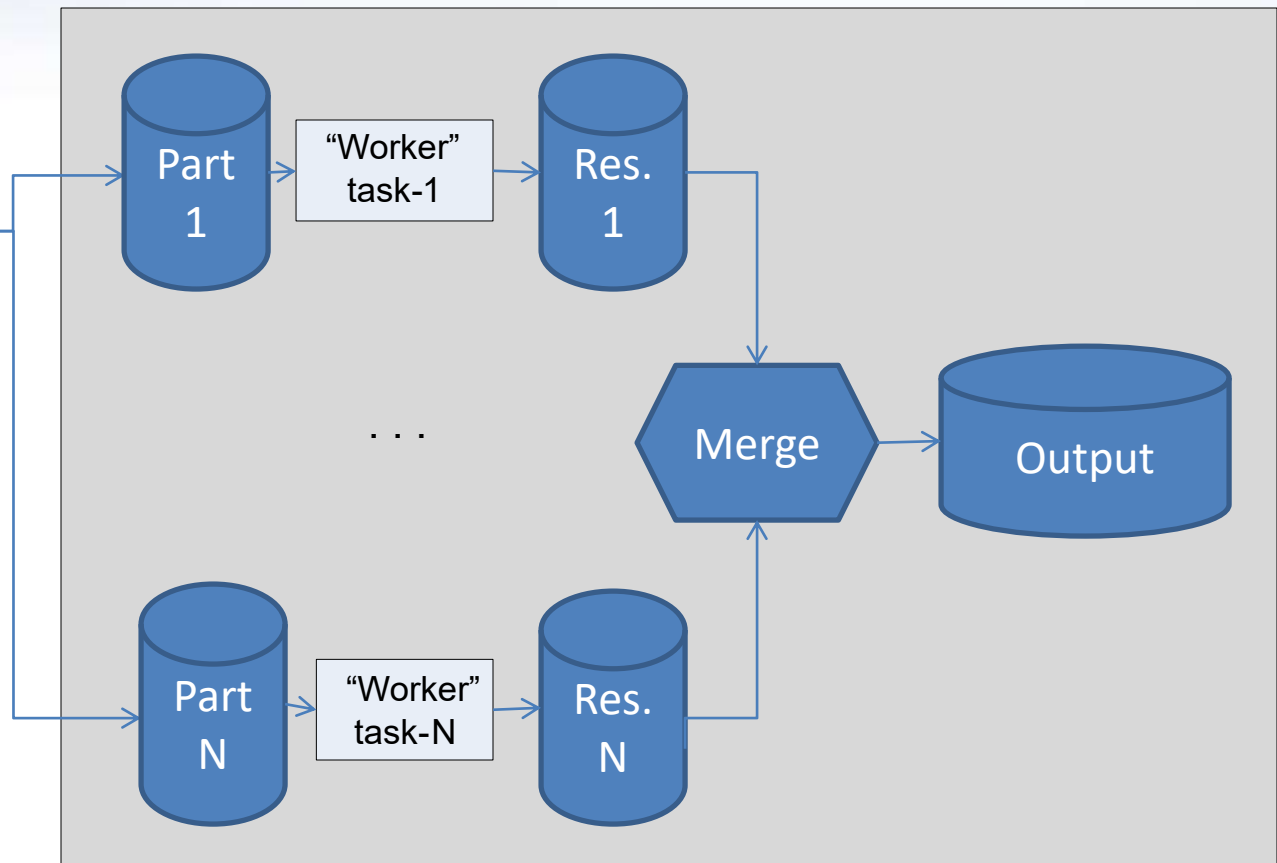


Data Scaling

Before scaling:
Time needed T



After scaling: Time needed $\sim T/N$



Scaling Techniques

Scaling Technique	Advantages	Disadvantages
Multi-threading, OpenMP	<ul style="list-style-type: none">• Multiple parallel threads within a node	<ul style="list-style-type: none">• Scaling is limited to cores on one node
MPI	<ul style="list-style-type: none">• Multiple parallel threads across one or more nodes	<ul style="list-style-type: none">• Overhead for I/O coordination and load balancing• In practice, all requested resources must be available to start• No checkpointing• Cannot exceed max capacity of the cluster
Scientific workflows, MapReduce, Spark, Hadoop	<ul style="list-style-type: none">• Scalable computational or data manipulation tasks on one or more nodes	<ul style="list-style-type: none">• Does not offer integrated approach for scaling multi-level nested loops or random number generation
Single loop parallelization	<ul style="list-style-type: none">• Multiple parallel tasks on one or more nodes	<ul style="list-style-type: none">• Does not parallelize multilevel nested loops
Array-based parallelization	<ul style="list-style-type: none">• Multiple parallel tasks across one or more nodes	<ul style="list-style-type: none">• Overhead for setup/convergence phases

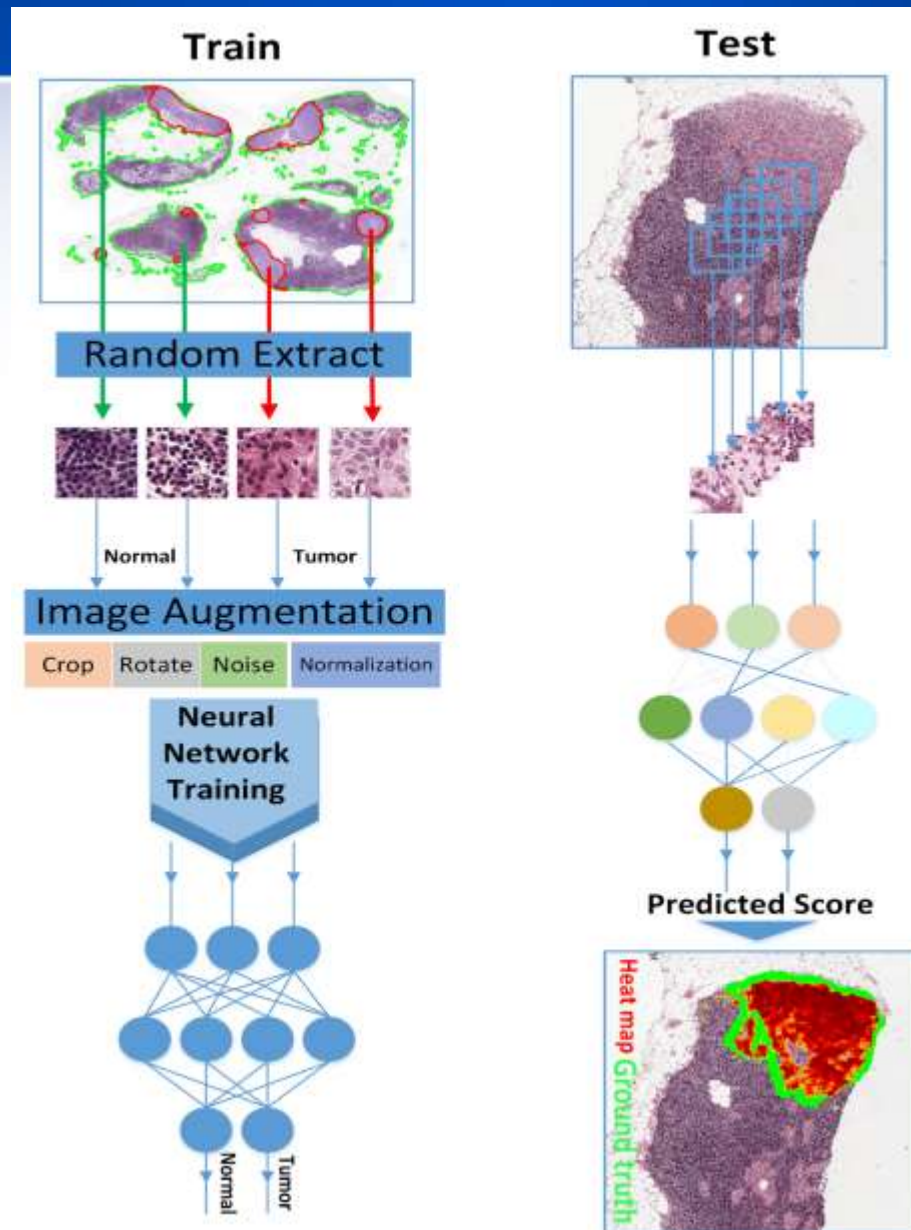
HPC DLNN Project



Weizhe Li, Weijie Chen, Mike Mikailov
CDRH/OSEL/DIDSR

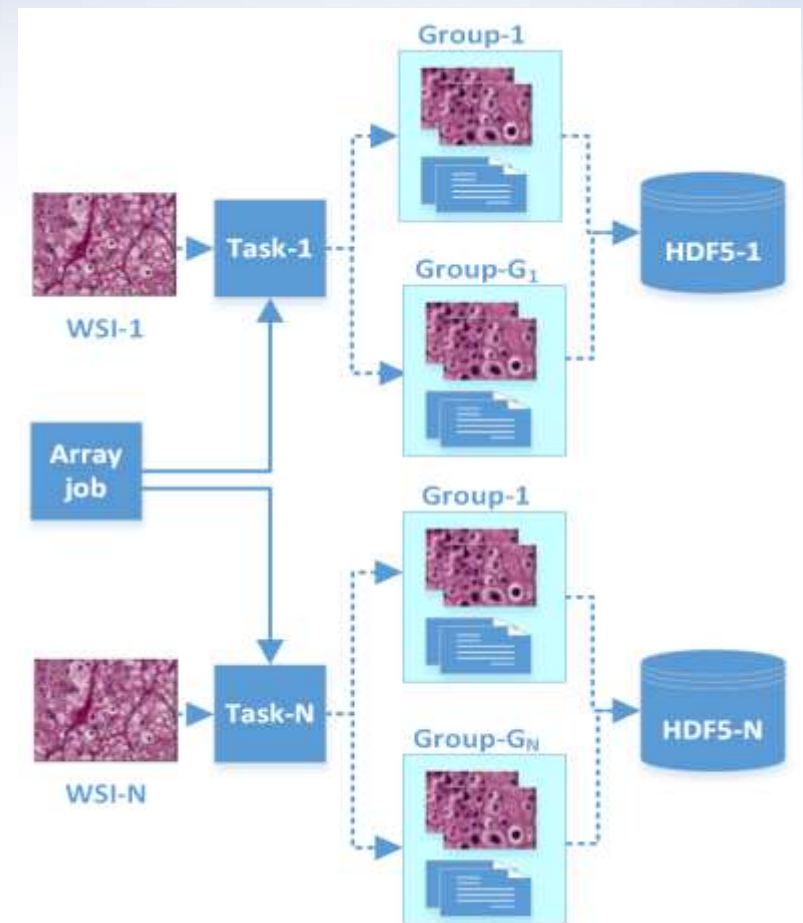
DLNN AI/ML with WSI data

- DLNN pipeline for digital pathology
- Training
 - Normal/tumor image patches randomly extracted
 - Normal (green)
 - Tumor (red)
 - Patches used to optimize NN using GPUs
 - HPC implementation generates pixel-wise heatmap via a sliding window



DLNN AI/ML with WSI data

- Scaling DDLN testing on HPC
 - Reformat/group patches into HDF5 file format for improved parallel I/O
 - Job splitting/scaling for parallel HPC implementations



DLNN AI/ML with WSI data

