

Storage Architectures, AI, and Biotech

RESEARCH

Multi-Vendor Webinar Tuesday March 23, 2021

Webinar Agenda



How Did Storage Architectures Perform for Biotech AI Modeling & What Can We Learn From This?

- **9:00-9:05** Ground Rules and Webinar Topic Introduction (G2M Research)
- **9:06-9:37** Sponsoring Vendor presentations on topic (8 minute each)
- **9:38-9:44** Panel Discussion Question #1
- **9:45-9:46** Audience Survey #1
- **9:46-9:52** Panel Discussion Question #2
- **9:53-9:54** Audience Survey #2
- **9:54-10:00** Panel Discussion Question #3
- **10:01-10:08** Audience Q&A (8 minutes)
- **10:09-10:10** Wrap-Up

G2M Research Introduction and Ground Rules

Mike Heumann (Managing Partner, G2M Research)

Artificial Intelligence and Biotech

- Artificial intelligence (AI) has demonstrated its value in a number of biotech areas
 - COVID screening, gene analysis
 - Understanding protein binding
 - Disease identification and diagnostics
- Al has enabled biotech companies to speed up drug analysis, design and screening
 - Screening potential drugs against thousands of molecules
 - Assessing impact of multi-drug treatments for complex diseases
 - Development of biomarkers





Storage Architectures for Biotech Al

- Most storage vendors are now optimizing their architectures for AI workloads
- Both cloud and on-prem solutions are now available for AI workloads
- These systems must provide large, scalable storage with high performance
- These storage architectures must also be able to provide data management to store training data sets and training results









Esteban Rubens Principal, Healthcare AI Practice <u>www.netapp.com</u>



WEKA

Greg Mazzu Sr. Systems Engineer <u>www.weka.io</u>





Scott Shadley VP of Marketing <u>www.ngdsystems.com</u>







Adam Marko Director, Life Sciences www.panasas.com

Andrew Bartko Executive Director, CMI/UCSD



Esteban Rubens Principal, Healthcare Al Practice <u>www.netapp.com</u>

Industry pain point

Data Scientists are doing very little Data Science



Data Scientist Work Distribution



Configuring hardware/platforms

Mass scale experimentation

Resource scheduling and assignments Kubernetes and Container management Open-source tools, plug-ins and dashboards Datasets import and management Collaboration and sharing

- Models and repo management Version control (models, data, git, etc.)
 - •Deployment

Production monitoring Automate Continual Learning







Bring a cloud consumption model to data science



- No need to reinvent the wheel: a turnkey solution for data science
- Minimize the time data scientists and engineers spend wrangling data
- Make tools available to data scientists while reducing the burden on IT
- Maximize resource utilization
- Avoid shadow Al
- Improve training data availability with caching



- Tools for data scientists built by data scientists
 - NetApp AI Control Plane: a full-stack AI data and experiment management solution
 - NetApp Data Science Toolkit: a Python program that makes it simple for data scientists and data engineers to perform advanced data management tasks
- Simplify results reproducibility with the Machine Learning Version Control framework
- Data Fabric: seamlessly move data
 - To and from clouds
 - Between different clouds
 - Between ONTAP systems (edge, core, and cloud)
 - ONTAP on AFF for high-performance model training paired with an S3 data lake
- Cloud cost optimization



Data Fabric: The right data in the right place at the right time, seamlessly



NetApp's portfolio spans the data pipeline





NetApp





PANASAS

Andrew Bartko Executive Director, CMI/UCSD Adam Marko Director of Life Sciences Panasas

www.panasas.com



Combating COVID-19 : Reducing SARS-CoV-2 transmission with early detection through wastewater monitoring

Dr. Andrew Bartko Executive Director of the Center for Microbiome Innovation (CMI) Professor of Practice – Bioengineering UC San Diego Observed and theoretical time lags between infection and detection of increasing SARS-CoV-2 transmission in wastewater and the health system











UC San Diego





Wastewater Monitoring Dashboard

Wastewater collection time period: 11AM PST February 12, 2021 -11AM PST February 13, 2021



Esri Community Maps Contributors, SanGIS, Esri, HERE, Garmin, SafeGraph, INCREMENT P, METI/NASA, ... Powered by Esri





Acknowledgements

Knight Lab

Expedited COVID IdenTification Environment (EXCITE Lab) at

UCSD

UCSD RTL Team

https://returntolearn.ucsd.edu/about/program-leadership/

Facilities Management

Point Loma WWTP

Center for Microbiome Innovation







PanFS[®] on ActiveStor[®] Ultra for Life Sciences HPC Storage G2M 2021

Adam Marko, Director of Life Science Solutions amarko@panasas.com

Researchers needing HPC

Source: NIH Biowulf Cluster Public Data

HPC users are growing at a faster rate than headcount

- HPC is needed in more research areas than ever before



©2020 Panasas, Inc.

What is Panasas?

PANASAS

- High Performance, Mixed Storage Media, File System Appliance
- NVMe, SSDs, and HDDs optimized without user config
- Scalable performance and data protection
- Data is auto balanced to ideal media transparently
- Users get seamless total performance from all storage media simultaneously

What makes Panasas Unique?



Dynamic Data Acceleration Enables Simultaneous Storage Access Data Types are Mapped to Ideal Storage Media



Thank You!

Adam Marko, Director of Life Science Solutions amarko@panasas.com info@panasas.com

APPENDIX

27

What is Panasas?

PANASAS

- A High Performance, Scalable, Parallel File System
 - Designed to meet dynamic workflows
 - Data protection, reliability, and performance increase as system grows
- In business for 20 years serving the HPC community
- Provide short- and long-term storage for HPC customers globally
- Greatly improved price/performance with Ultra release

Research Cost of Ownership (RCO)



The Effect of Infrastructure on Research Productivity

Researchers

- Publications and research outcomes
- Timetables and collaborations
- Poor morale and frustration

IT

- Annual goals
- Daily responsibilities
- Impact on other projects

Leadership

- Cost of reduced staff productivity
- •Effect on research goals
- Reputation

Using RCO

- Ask your Researchers, IT, and Leadership questions
- Don't make decisions based on side-by-side cost comparison alone



Panasas for AI and ML

Al without all flash

Storage System (Half Rack)	V100 GPUs	Capacity	Cost
Panasas ActiveStor® Ultra	42	1.5PB raw	\$
Dell Isilon F800	52	768TB raw	\$\$\$\$

- PanFS supports 100s of GPUs per rack for demanding AI/ML applications
 - Configured for AlexNet with TensorFlow using ImageNet dataset
- Compared to Isilon F800 All-Flash, PanFS provides:
 - ~ 2x the capacity
 - 1/5th the cost
 - Similar number of GPUs supported





https://www.delltechnologies.com/resources/en-us/asset/white-papers/products/storage/h17361_wp_deep_learning_and_dell_emc_isilon.pdf

Greg Mazzu Senior Systems Engineer <u>www.weka.io</u>

1000

WEKA

Welcome to Weka

Our Mission

Make storage a utility by delivering simplicity, speed, scale, and better economics



WEKA

8 of the Fortune 50 are customers

Backed By Industry Leaders





Limitless Data Platform for Health & Life Sciences



Boston Pharmaceutical Solution



Weka Deployment at Scale - Genomics England



Weka on AWS for NVIDIA Clara Parabricks

📀 nvidia. Develo	PER					Q Join Login
NVIDIA Developer I	Blog		e d	EVELOPER NEWS		🛩 FOLLOW US
AI / DEEP LEARNING	AUTONOMOUS MACHINES	AUTONOMOUS VEHICLES	DATA SCIENCE	GRAPHICS / SIMULATION	HPC IVA/IOT	NETWORKING

AI / DEEP LEARNING

Analyzing Genome Sequence Data on AWS with WekaFS and NVIDIA Clara Parabricks Pipelines

By Robert Clark and Bob Bakh | October 2, 2020 * Tags: AWS, Clara, Computational Chemistry, genomics, Healthcare, miligate, Parabricks, Weka 🥏 2 Comments

Whole genome sequencing has become an important and foundational part of genomic research, enabling researchers to identify genetic signatures associated with diseases, differentiate sequencing errors from biological signals, and better characterize the genomes of various organisms. With the ongoing COVID-19 pandemic threatening the globe, characterizing, and understanding genomes is now more crucial than ever. Commercially available, next-generation sequencing platforms allow researchers to decode an entire human genome in less than a day. This helps in understanding the susceptibility with infection by SARS-CoV-2; can be used as the basis for vaccine creation; and can be used for therapy selection for an individual based upon their unique genetic signatures, along with many other use cases.

Traditionally, sequencing a whole human genome takes multiple days coupled to a heavy compute power using CPU resources. The <u>Genome Analysis Toolkit (GATK)</u>, developed by the Broad Institute, is a multi-purpose software suite for analyzing DNA_ and RNA-based sequence data with the primary goal to identify genetic variants. It is generally considered the industry standard toolset. The GATK suite contains several tools, including the Base Quality Score Recalibration (BQSR), Burrows Wheeler Aligner Maximal Effect Match (BWA MEM) for aligning and calibrating genomes, and HaplotypeCaller, which efficiently identifies variants in sequences. These tools are all run as sub-stages in the germline pipeline that identifies germline variants and is one of the algorithms used in this post.

NVIDIA Clara Parabrick Pipelines replicate the functionality of GATK while harnessing the power of GPUs to provide the fastest method to analyze genomes. The software, leveraging the NVIDIA CUDA libraries to accelerate key algorithms, dramatically reduces the time required to analyze a sequence compared to the CPU-only GATK tools and also includes the <u>DeepVanant algorithm</u> developed by researchers at Google. Combined with the power of <u>Amazon Web Services</u> [AWS]



WEKA

and NVIDIA storage partner WekalO (Weka), the high-performance sequence analysis of NVIDIA Clara Parabricks integrates with the convenience of the cloud to offer a robust, fast, and simple solution to enable faster genomic research over existing state-of-the-art CPU-based solutions. It offers over 33x improvements in performance with over 99.9% concordance of results.





Scott Shadley VP of Marketing www.ngdsystems.com

G NGD systems



Computational Storage NVMe SSDs with Compute on Drive

Driving Compute and Storage in any Environment



The Market Needs a New Way to Look at Storage.



Pain PointsPhysical SpaceAvailable PowerScaling MismatchBottleneck Shuffle

Traditional storage architectures are in trouble.

Current Compute Systems are limited to a few CPU count per system The Processing Unit (CPU/GPU/TPU) all require scale for BioTech



Finding ways to process data more efficiently is needed This support needs to come **without wholesale change**

39

The Market Solution Using Computational Storage.



Value Add Distributed Processing **Faster Results** Lower Power Smaller Footprint Scaling compute resources within storage provides accelerated results

Computational Storage resources 'offload' work from the limited CPU count

Seamless architectures create new 'servers' for more effective analytics



This 'Server in a Server' Architecture provides value across many use cases

Additional CPU resources for the cost of Storage without added Resources

ARM NN Deployment for Parallel & Distributed Processing





- More than **100k data points per drive** are processed without sending the data to host. Saving Time, Network and Power.
- The output of the application is written directly to a MongoDB database implemented inside In-situ engine, and user can access the output data using MongoDB APIs.



VIDEO DEMO WITH Accelerated response

41

NGD NVMe SSD Products at a Glance.

- Large breadth of SSD solutions and capacity options
- Leading TB/W Energy Efficiency
- Industry's only **16-Channel** 14nm SSD SoC
- Industry's Largest capacity NVMe SSDs
- Quad-Core Computational Storage CPUs

Form Factor	Availability	Raw Capacity TLC (TB)	MAX Power (W)
M.2 22110	NOW	up to 8	8
U.2 15mm	NOW	up to 32	12
EDSFF E1.S	NOW	up to 12	12







42

Computational Storage, Some Real World Results.



NGD systems

Al Inference Offload

- Similarity search, Tracking

Distributed Machine Learning

- Identification, Training, Tagging

Data Search, Compare

- Finding the Needle, Making sure it is!

Database Acceleration

- Data is stored somewhere, drive it efficiently

Protein Sequencing – BLAST® Accelerated



• DNA and Protein alignment Database Management

EXACTION The Basic Local Alignment Search Tool (BLAST) finds regions of similarity between sequences.

- Not balanced dataset
 - o Computation time varies a lot across files
 - o Different number of sequences per file
- By combining with Computational Storage SSDs and using the 4 cores per drive, you gain up to 100% more performance at no cost in CPU or Memory



Accelerated ML, Faster, More Efficient



• Four neural networks Evaluated

- o MobilenetV2
- o NASNet
- o SqueezeNet
- InceptionV3Quad-core
- Training data stored on CSDs
- Tested with 24 CSDs
 - \circ 32TB capacity each
 - o Quad-core CPUs on Drive
 - o 4x NEON SIMD engines

• Using an AIC 2U-FB201-LX server

Intel[®] Xeon[®] Silver 4108 CPU
32GB DRAM



45

Panel Questions and Audience Surveys

Panel Question



- How does storage for biotech AI applications differ from other AI applications and/or other HPC applications?
 - NetApp
 - Panasas
 - Weka
 - NGD Systems

Audience Survey Question #1



 We use AI and ML widely in our business/organization: 	20%
 We have several problem domains which utilize AI and ML: 	20%
 We have a few problem domains in which we utilize AI and ML: 	0%
 We are exploring the use of AI and ML for new workloads, but have not yet implemented any AI/ML solutions: 	30%
 We have no plans to utilize AI and ML workloads in our organization in the near future: 	30%
 Don't know: 	

Panel Question #2



- We have all heard about the importance of managing training data sets, even as they evolve over time. What technologies are key to integrating and managing different data types?
 - Panasas
 - Weka
 - NGD Systems
 - NetApp

Audience Survey Question #2

 When looking at implementing AI and ML storage solutions, who do you look to for guidance on hardware and software choices (check all that apply):

 My current IT server vendor(s): 	21%
 My current IT storage vendor(s): 	43%
 My current IT integrators: 	0%
 AI/ML application/framework vendors: 	43%
 AI/ML application/framework integrators: 	14%
Other:	21%

Panel Question #3



- All of your companies provide on-premise storage solutions. How does the cloud fit in to storage architectures and solutions for Al applications?
 - Panasas
 - Weka
 - NGD Systems
 - NetApp

Audience Q&A



Thank You For Attending!