

Who am I?

I'm the co-founder and CEO of **Lambda**. We've built large scale GPU clusters for the Fortune 500, the world's leading academic research institutions, and the DOD.

I'm also the lead architect of the **Lambda Echelon**, a turn-key GPU cluster. This talk is based on the Lambda Echelon reference design whitepaper and the experience we've gained deploying these large clusters.



Overview

- Introduction
- Cluster use case overview
- Cluster design
 - o Compute
 - Storage
 - Cluster Networking
 - o Data center floor planning & power distribution checklist
 - Software
- Rack design
 - o Rack elevations
 - Power basics
 - o Rack power distribution
- Node design
 - GPU selection
 - Node PCIe topology / NVLink topology / real life examples
 - o GPUDirect RDMA
- Putting it all together
- Citations



5 Stages of GPU Cloud Grief



It all starts with the **shock** of an expensive cloud bill.

Greetings from Public Cloud Provider,

This e-mail confirms that your latest billing statement, for the account ending in 1234, is available on the Cloud web site. Your account will be charged the following:

Total: \$18,445.92

You can see a complete break down of all charges on the Billing & Cost Management page.

Thank you for using the public cloud.

Sincerely, Public Cloud Provider



Stage 1 - Denial

"This won't happen again next month."



Stage 2 - Anger "The bill doubled again!"



Stage 3 - Bargaining with your account manager.

Hi Team,

Are you available in the next few days for a 15 min call to discuss lowing your public cloud spend using reserved instances?

Based on your January usage, purchasing a 1-year reserved instance would save you up to 20% on your public cloud bill!

Thank you for being a loyal public cloud customer.

Cheers,
Public Cloud Provider



Stage 4 - Depression

"Spot instances and reserved instances aren't enough, this is hopeless."



Stage 5 - Acceptance

"GPU cloud services are expensive. Managing hardware is scary."

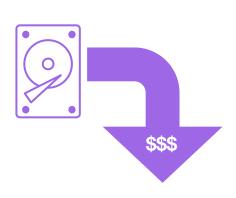


Solution

Learn how to build a GPU cluster.



Good reasons to consider building an on-prem cluster



Big data sets & expensive egress



Data sovereignty & security



More compute for less money



Rule Number 1:

Know your ML team.



Your use case should inform the cluster design

14

Hyperparameter search

o "Finding the best model."

Large scale distributed training

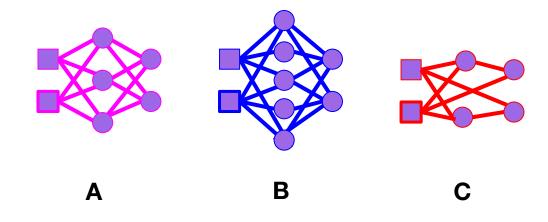
"Quickly training a model."

Production inference

"Deploying the model at scale to production."

Hyperparameter search

Which neural network performs best (highest accuracy)?

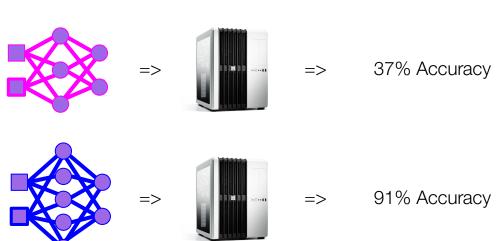


Hyperparameter Search

Answer: _____?____

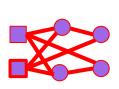


Hyperparameter search scoring



Hyperparameter Search

lambdalabs.com



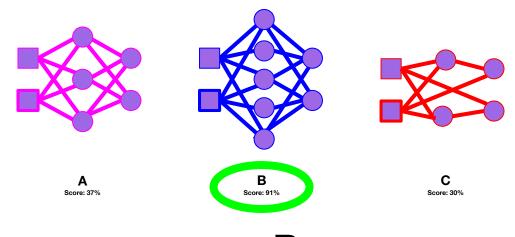


=> 30% Accuracy



Hyperparameter search

Give a score to each, select the highest score.



Hyperparameter Search

Answer:



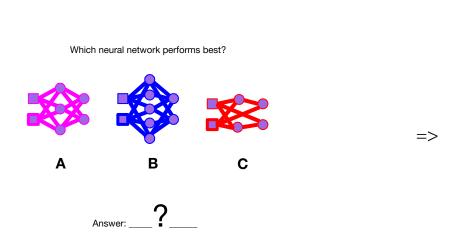
Each training run can take 1 to 48+ hours!



Hyperparameter Search



Most people start small

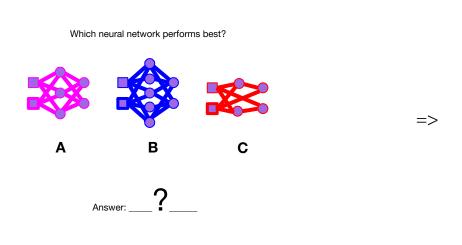




Hyperparameter Search



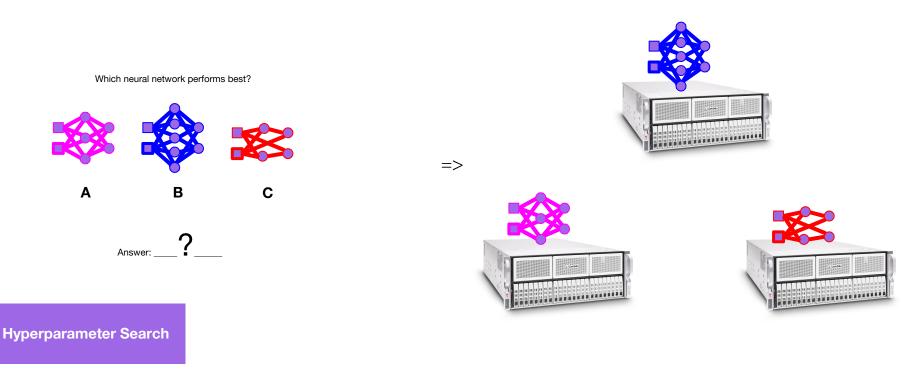
Then scale to faster GPUs





Hyperparameter Search

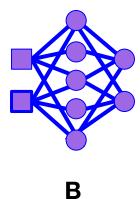
Then scale to multiple servers to run the jobs simultaneously





Large scale distributed training

Train this model as fast as possible

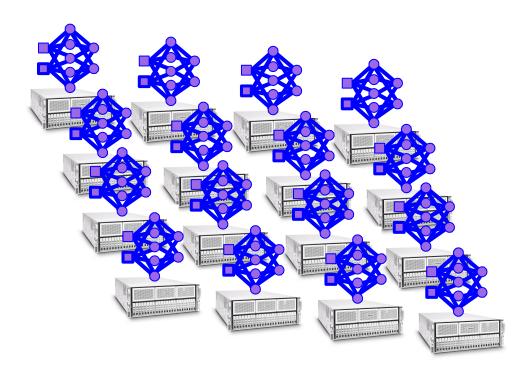


22

Large Scale Distributed Training



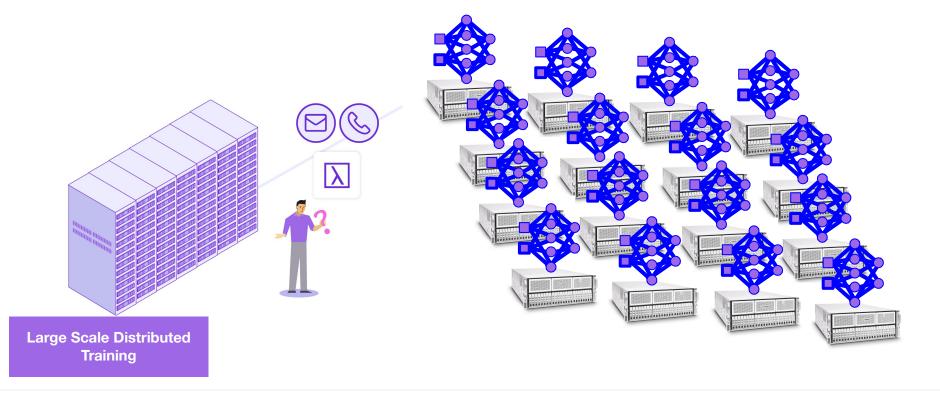
Have multiple servers help with the training



Large Scale Distributed Training

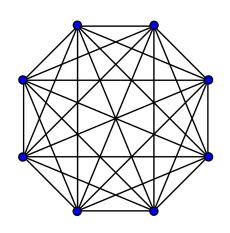


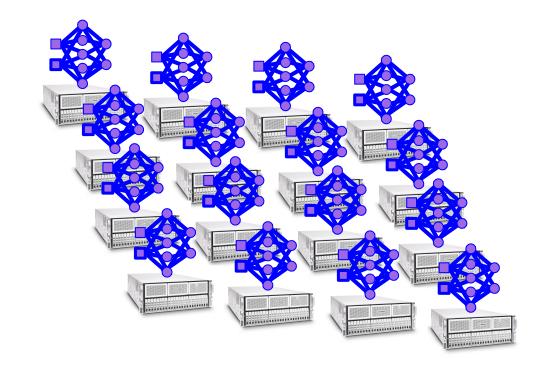
Coordinating the servers is hard





Lots of server-to-server communication = high speed network required

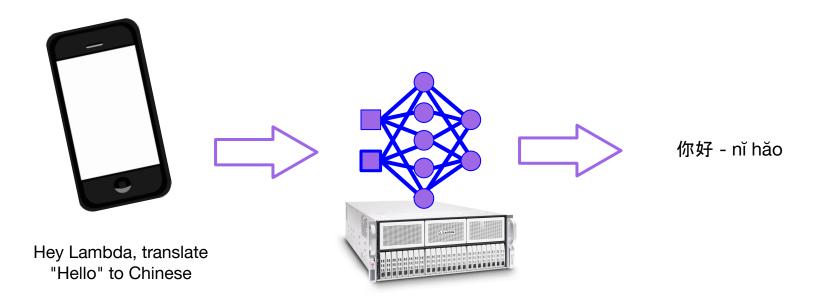




Large Scale Distributed Training



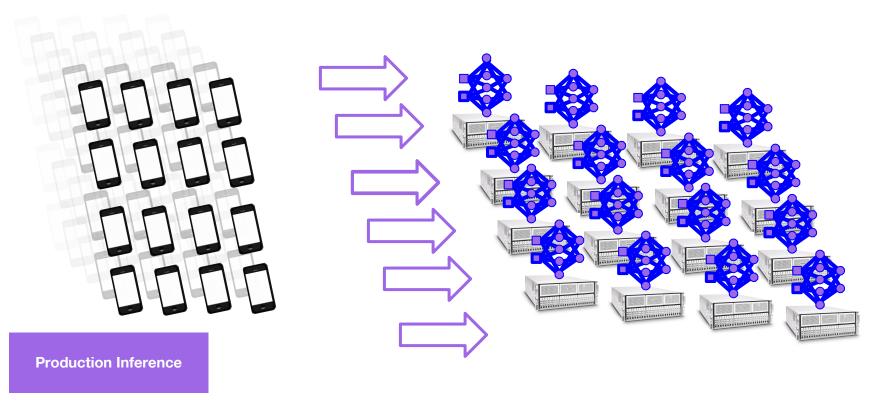
Production inference



Production Inference

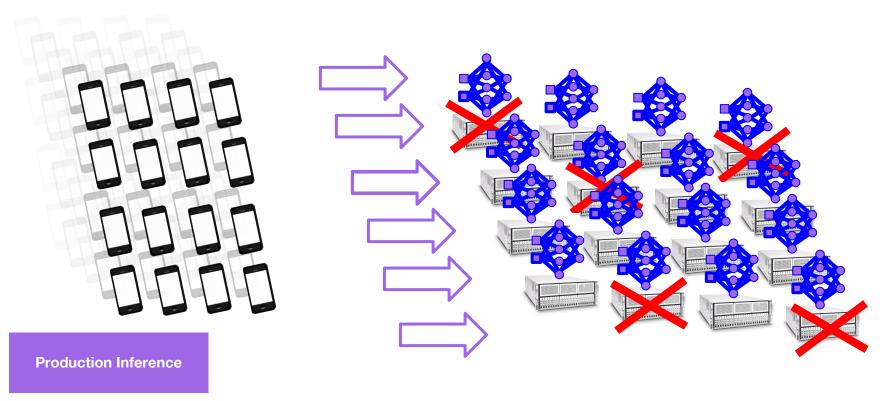


Inference clusters handle thousands of simultaneous requests





Inference clusters must be robust to outages



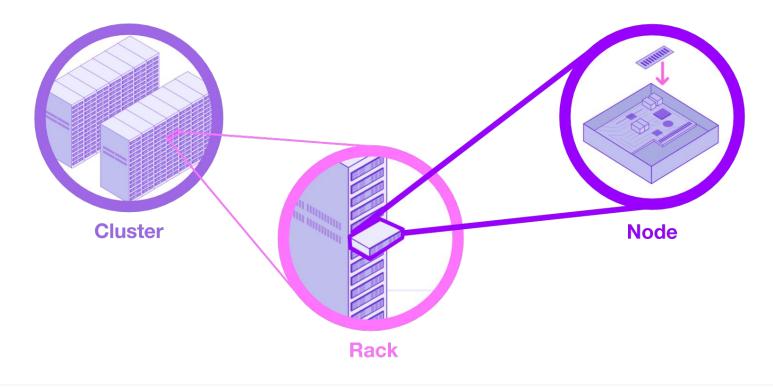


What's the right cluster for each use case?

	Hyperparameter search	Large scale distributed training
Node to Node Bandwidth	Low	High
Key Metric	Training throughput / \$	Time to train a single large model
Price	\$\$	\$\$\$\$
1+1 Redundancy	Optional	Optional
Operating Mode	Usually offline, job-queue based	Usually offline, job-queue based
	Production inference	
Node to Node Bandwidth	Low	
Key Metric	High availability & throughput	
Price	\$\$\$\$\$	
1+1 Redundancy	Critical	
Operating Mode	Often online, real time results	



Three levels of abstraction





The three levels of cluster design

- 1. **Cluster design:** the highest level of abstraction. Entire racks are just dots on the screen. Data center floor plans, capacity planning (power, network, storage, compute), and network topologies are the main output products at this level.
- 2. **Rack design:** Rack elevations describe the layout and exact position of individual nodes in a particular rack. Cable and port counts are important details at this layer.
- 3. **Node design:** Individual node bills of materials and cluster design goals drive component selection for each node.

Software: it's important to remember that there is a two way design dependency between the cluster hardware and the software that runs on it. Selecting the right software depends on the use case and scale of your planned deployment.



Work products for each layer

Abstraction Layer	Work Products	Goals
Cluster	 Cluster BOM (Bill of Materials) Data center floor plan & checklist Capacity plan Network topology Installation plan 	 Working within constraints of data center floor space and power availability Connecting all of the racks together Determining the overall storage/compute/network capacity
Rack	Rack BOMRack elevation	 Properly distributing the capacities planned for in the cluster stage Rack level power density & capacity Node placement within rack
Node	Node BOM	Hitting the performance requirements defined by the compute capacity of the cluster
Software	Software architectureSoftware installation plan	Tying all of the hardware together so that it's easy to use and administer

32







Cluster architecture

Cluster architectures have five components:

- Compute: to do the actual work. GPU & CPU nodes.
- 2. **Storage:** to serve data sets and store trained models / checkpoints.
- 3. **Networking:** multiple networks for compute, storage, in-band management, out-of-band management.
- 4. **Data center power distribution & floor planning:** understanding the electrical and physical layout of the deployment data center informs the rack elevations, cable lengths, and networking.
- 5. **Software:** cluster orchestration, job scheduling, resource allocation, container orchestration, and node level software stacks.



Compute

35

You need to answer three questions:

- 1. How much compute do you need?
- 2. How fast do you want the compute to talk to each other?
- 3. How fast do you need the compute to be able to talk to the storage?

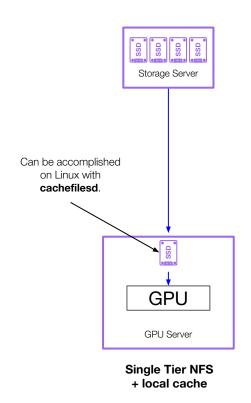


Storage

- Often, storage clusters will become the bottleneck in a highly optimized cluster. It's important to remove this bottleneck because it reduces the utilization of the most expensive part of the cluster, the compute.
- There are a two ways to design a storage cluster:
 - Work with a storage partner
 - Roll your own

36

Storage architecture diagrams



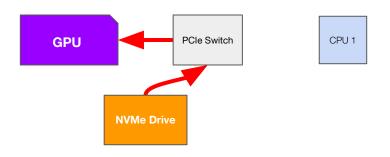
Storag Metadata Server Storage Server **GPU GPU Server** Parallel cluster file system + local cache

HDD Server SSD Server **GPU GPU Server**

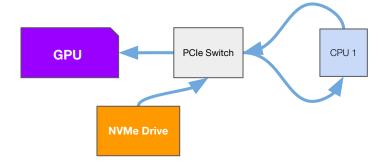
Tiered storage system + local cache



GPUDirect Storage



Data pathway with GPUDirect Storage (Directly to the GPU via PCle switch)



Data pathway without GPUDirect Storage (Additional copy to CPU memory)



Storage choices

Open source storage options











Proprietary storage options



















Cluster networking



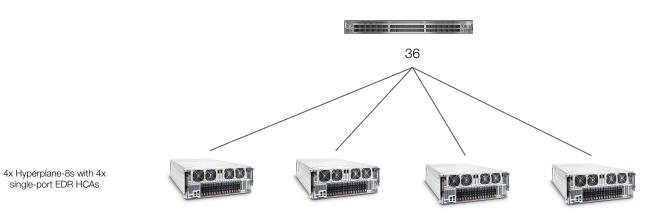
Networking basics: switches

MSB7800 - 36 port EDR (100Gb/s) IB switch - 7.2Tb/s aggregate switch throughput 7.2 Tb/s / 36 ports = 7200 Gb/s / 36 ports = **200 Gb/s / port => (100 Gb/s EDR IB)**



Networking basics: single switch topology

Non-blocking 100% port to port speed



4*4 = 16 total IB cables

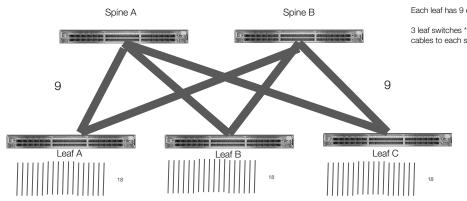


Networking basics: spine and leaf topology

Non-blocking fat tree topology

To guarantee non blocking, ensure equal number of "southbound" and "northbound" ports on the leaf switches

36/2 = 18 southbound and 18 northbound



12x Hyperplane-8s with 4x single-port EDR HCAs



4*12 = 48 total IB cables

Each leaf has 9 cables going northbound, 3 leaf switches

3 leaf switches * 9 northbound cables = 27 northbound leaf switch cables to each spine switch.

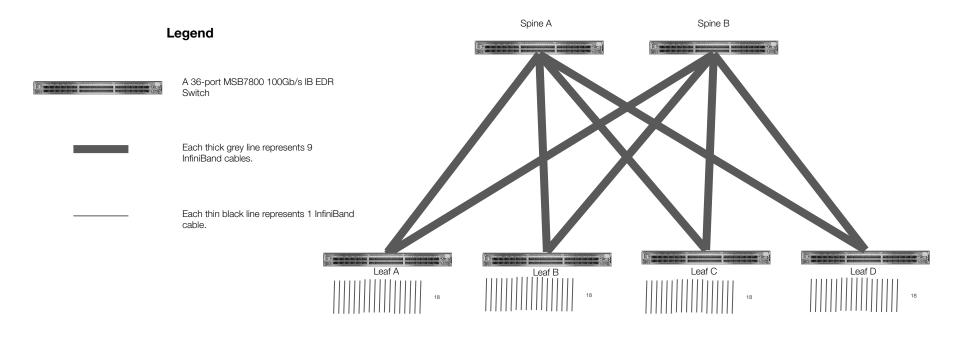
48 needed cables / 18 southbound = 2.666, round up to 3 => 3 leaf switches to provide full 48 ports.

54 total ports - we're using 48/54.

Total southbound bandwidth of 5.4Tb/s available.



100% port-to-port bandwidth spine & leaf topology





Networking basics: director switches

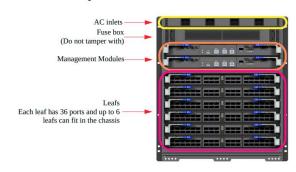
45

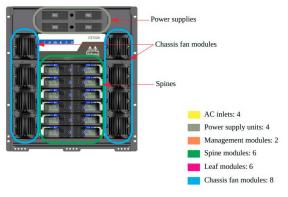
- Use a backplane instead of cables to build out the internal (spine <> leaf) network topology.
- To get 100% non-blocking bandwidth in the CS7520 you NEED to have all 6 spines installed.
- CS7520 has 216 southbound ports (and thus 216 northbound ports).
- Assuming each spine has only 36 ports, how many spines do you need to support 216 ports?
 - 36 ports / spine, 216 ports needed
 - 216 ports / 36 (ports / spine) = 6 spine switches

For more info see:

https://www.mellanox.com/related-docs/prod_ib_switch_systems/C S7520 Dismantling Guide.pdf

Figure 1: Front and Rear View of the CS7520

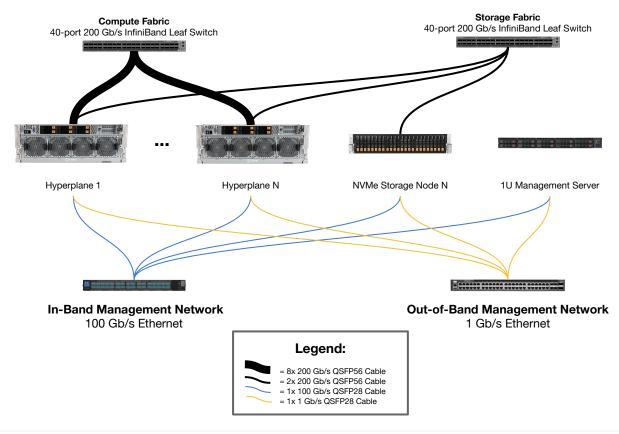






Lambda Echelon Network Topology

(Single Rack Configuration)

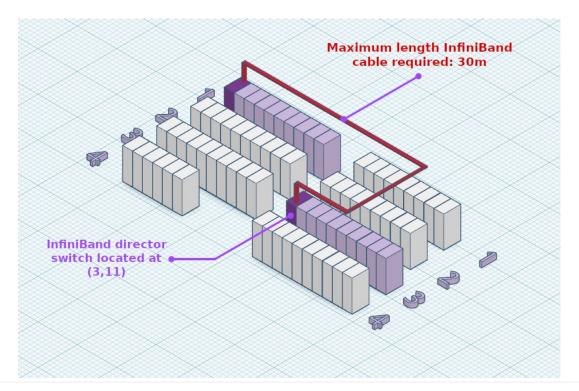




Data center floor planning & power distribution

Checklist:

- What rack does your DC provide?
- Max power (kW) per rack?
- What PDUs are provided, if any?
- What's the floor plan of the DC?
- Pull floor plan into DCIM/CAD to calculate cable lengths.



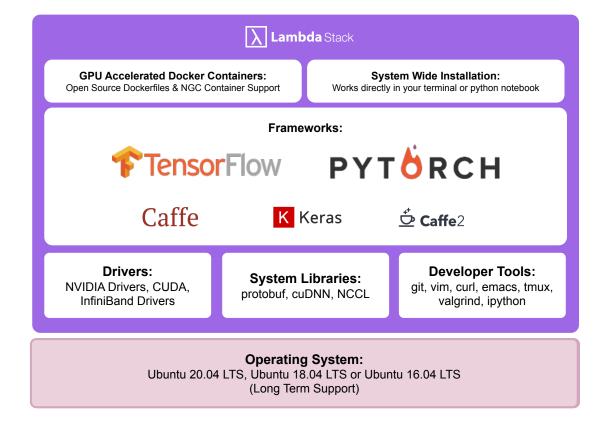


Software

There's a lot of software you'll want to run on your cluster. Some of this software is installed on the compute nodes, some of it is installed on the management nodes. Depending on your use case, these are some common packages you'll encounter.

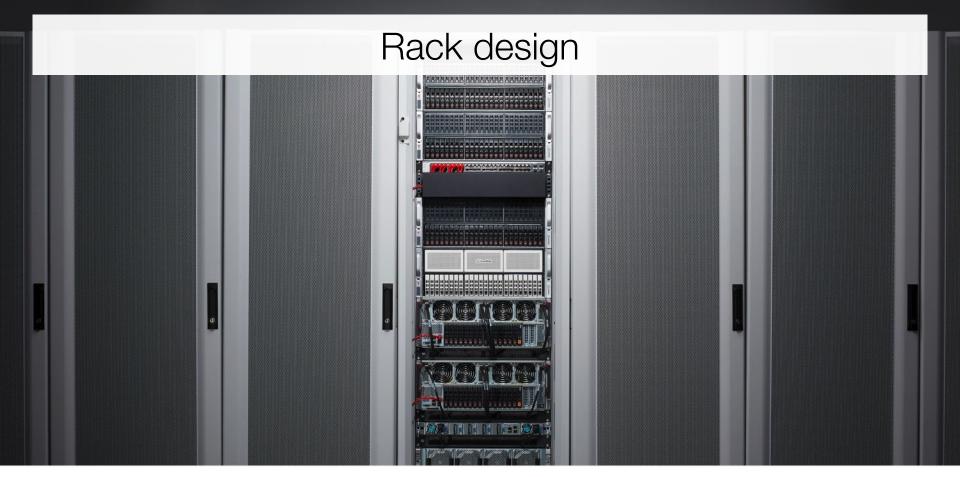
- Experiment management software (Weights & Biases, Determined.ai, Spell.ml, etc.) manage experiments and version control
- Notebooks Jupyter Lab or iPython notebooks
- Containers Docker mostly but also Singularity
- Kubernetes manage instances and execute containers on them
- Job scheduler (SLURM, Kubeflow)
- Logging (Prometheus, Grafana, etc.)
- **Production inference** (TensorFlow Serving, TensorRT)
- Lambda Stack drivers, CUDA, TensorFlow, upgrade management, Dockerfile templates, etc.



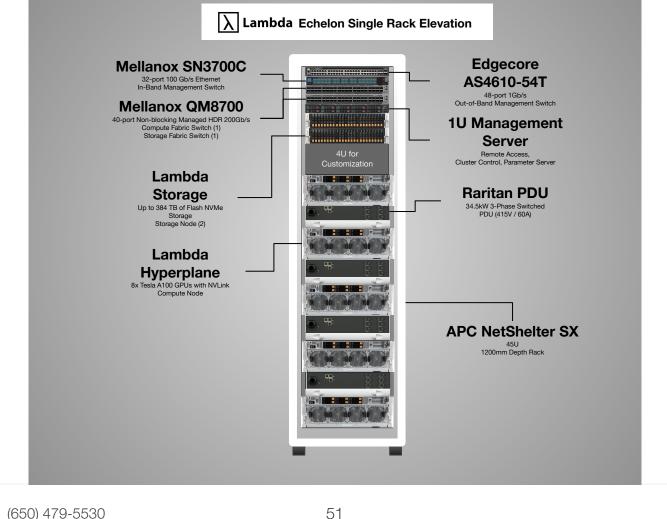


For more info on Lambda Stack, see: https://lambdalabs.com/lambda-stack-deep-learning-software

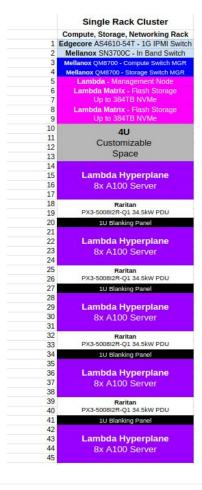














Lambda Echelon - Single Rack Cluster - Bill of Materials

Description	Qty	TDP Draw (kW)	Ext. TDP (kW)	
Lambda Hyperplane - 8x A100 Server	5	6.5	32.5	
1U Management Node	1	1	1	
2U Flash Storage Server	2	0.75	1.5	
Mellanox SN3700C - 32 port - 100 Gb/s Ethernet - In-Band Switch	1	0.25	0.25	
Mellanox QM8700 - 40 port - 200 Gb/s IB- Compute Fabric Switch	1	0.8	0.8	
Mellanox QM8700 - 40 port - 200 Gb/s IB - Storage Fabric Switch	1	0.8	0.8	
Edgecore AS4610-54T - 48 port - 1 Gb/s Ethernet - IPMI Switch	1	0.1	0.1	
Mellanox MFS1S00-H003E - 3m active 200 Gb/s IB cables	49	0	0	
Mellanox MFA1A00-C003 - 3m active 100 Gb/s Ethernet cables	8	0	0	
Cat 5e Ethernet Cables - 1m/2m/3m	16	0	0	
Line A - Raritan PX3-5008I2R-Q1 - 34.5kW PDU 415V 60A Input IEC 60309 3P+N+E 6h 60A (4P5W)	2	0	0	
(Optional Line B) Raritan PX3-5008I2R-Q1 - 34.5kW PDU 415V 60A Input IEC 60309 3P+N+E 6h 60A (4P5W)	2	0	0	
1U Blanking Panel	4	0	0	
APC Netshelter SX 3305W - 45U - 600mm x 1200mm Rack	1	0	0	
C14 to C13 cables (2 ft)	14	0	0	
C13 to C20 cables (2 ft)	20	0	0	
Rack TDP Draw (kW):				



Lambda Echelon - Single Rack Cluster - Bill of Materials

Description	Qty	TDP Draw (kW)	Ext. TDP (kW)
Lambda Hyperplane - Bx neer most important 1U Management Node TO CONTROL OF THE PROPERTY OF	5	6.5	32.5
1U Management Node	1	1	1
2U Flash Storage Server	2	0.75	1.5
Mellanox SN370 C 32 phr - 100 cb/s Eper 4 - In-Gan Switch This page . Mellanox QM8700 - 40 port - 200 Gb/s IB- Compute Fabric Switch	1	0.25	0.25
	1	0.8	0.8
Mellanox QM8700 - 40 port 600 Gb/s IB - Stor ge Fabric Switch	1	0.8	0.8
Mellanox QM8700 - 40 port 680 Gb/s IB - Storge Fabric Switch Edgecore AS4610-54T - 48 pr 11 kb s arr t - IF Al arc C	1	0.1	0.1
Mellanox MFS1S00-H003E - 3m active 200 Gb/s IB cables	49	0	0
Mellanox MFA1A00-C003 - 3m active 100 Gb/s Ethernet cables	8	0	0
Cat 5e Ethernet Cables - 1m/2m/3m	16	0	0
Line A - Raritan PX3-5008I2R-Q1 - 34.5kW PDU 415V 60A Input IEC 60309 3P+N+E 6h 60A (4P5W)	2	0	0
(Optional Line B) Raritan PX3-5008I2R-Q1 - 34.5kW PDU 415V 60A Input IEC 60309 3P+N+E 6h 60A (4P5W)	2	0	0
1U Blanking Panel	1	0	0
APC Netshelter SX 3305W - 45U - 600mm x 1200mm Ra	1	0	0
C14 to C13 cables (2 ft)	14	0	0
C13 to C20 cables (2 ft)	20	0	9
		Rack TDP Draw (kW):	36.95

\(\) Lambda

Back to basics

Single phase systems

P (watts) = V (volts) * I (amps)

We use I from the French, intensité du courant.

Because that's what André-Marie Ampère used.

3-phase systems

This simplifies to P = sqrt(3) * V * I

Because 3 / sqrt(3) = sqrt(3)

P = 3 * V / sqrt(3) * I

Real life 3-phase systems

P = sqrt(3) * V * I * 0.8

55

It's very common to see an 80% regulatory derating factor applied to PDUs.





How do PDU manufacturers calculate power capacity?

56

From the APC8966 Data Sheet:



Example PDU:

APC 8966

Input frequency

50/60 Hz

Number of Power Cords

1

Load Capacity

17300VA

Maximum Input Current

60A

Maximum Line Current

48A

Regulatory Derated Input Current (North America)
48A

Iominal Output Volta

Nominal Output Voltage 208V

Nominal Input Voltage

208V 3PH

Input Connections

IEC 60309 60A 3P + PE

P = sqrt(3) * V * I * 0.8

Plug in the numbers from the data sheet:

P = sqrt(3) * 208 * 60 * 0.8

P = 17292.7953

P = 17.3kVA

See how they derate the maximum input current of 60A to the "Regulatory Derated Input Current (North America)" 48A? That's 48A = 60A * 0.8.

That's where the 0.8 came from in the previous slide.

Some common PDU input plugs



IEC60309 - 60A 3-phase plug - 208V Blue means the system is between 200 and 250V.



IEC60309 - 60A 3-phase plug - 415V Red means the system is above 400V.



NEMA L15-30P 30A 3-phase plug - 208V



Receptacles & Plugs

	•		S
Plug Photo	Name	Plugs into	Receptacle Photo

IEC C13 Plug

IEC C14 Plug

IEC C19 Plug

IEC C20 Plug

(650) 479-5530

lambdalabs.com

IEC C14 Receptacle

(on server)

IEC C13 Receptacle

(on PDU)

IEC C20 Receptacle

(on server)

IEC C19 Receptacle

(on PDU)

IEC stands for International Electrotechnical Commission, an international standards organization headquartered in Switzerland.

58

Max Amps

15A

(Max power ~2.5kW)

15A

(Max power ~2.5kW)

20A

(Max power ~3.2kW)

20A

(Max power ~3.2kW)

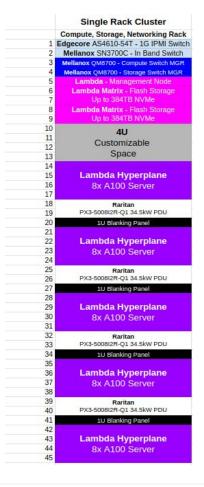
Lambda Echelon - Single Rack Cluster - Bill of Materials

Description	Qty	TDP Draw (kW)	Ext. TDP (kW)
Lambda Hyperplane - 8x A100 Server	5	6.5	32.5
1U Management Node	1	1	1
2U Flash Storage Server	2	0.75	1.5
Mellanox SN3700C - 32 port - 100 Gb/s Ethernet - In-Band Switch	1	0.25	0.25
Mellanox QM8700 - 40 port - 200 Gb/s IB- Compute Fabric Switch	1	0.8	0.8
Mellanox QM8700 - 40 port - 200 Gb/s IB - sharate Carrie Witch YOU UNCESSED STATE OF SWITCH YOU WINDOW YOU	tan	0.8	0.8
Edgecore AS4610-54T - 48 port - 1 Gb/s Etilernet - First Switch	tai 1	0.1	0.1
Mellanox MFS1S00-H003E - 3m active 200 Gb/s IB cables	49	0	0
Mellanox MFA1A00-C003 - 3m active 100 Copy Charles Turch Work	COF	0	0
Cat 5e Ethernet Cables - 1m/2m/3m	9016	0	0
Line A - Raritan PX3-5008I2R-Q1 - 34 5kW PDU 415V 60A Input IEC 60309 3P+N+E 6h 60A (4P5W) DT (Optional Line B) Raritan PX3-5008I2F-Q1 - 31 5kW PDU	oriz		0
(Optional Line B) Raritan PX3-500812F-Q1 - 3 L 3 L 3 L 3 L 3 L 3 L 3 L 3 L 3 L 3		113. ₀	0
1U Blanking Panel	1	0	0
APC Netshelter SX 3305W - 45U - 600mm x 1200mm Rack	1	0	0
C14 to C13 cables (2 ft)	14	0	0
C13 to C20 cables (2 ft)	20	0	0
Rack TDP Draw (kW):			



Lambda Echelon rack elevations





40 GPUs

Training time for Mask R-CNN on MSCOCO: 25 mins

Source: https://mlperf.org/training-results-0-7/





160 GPUs

Training time for Transformer WMT E-G Big: ~1.02 mins

Source:

https://mlperf.org/training-results-0-7/





800 GPUs

Training time for ResNet 50 on ImageNet: 1.06 mins

Source:

https://mlperf.org/training-results-0-7/ https://github.com/NVIDIA/DeepLearningExamples/blo b/2de29455a696b5608e24daa95f579f689b2d59e1/Pv Torch/Translation/Transformer/README.md



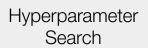
Node design





Choosing the right GPU depends on your use case







Large scale distributed training



Production inference



GPU benchmarking is hard!



https://lambdalabs.com/blog/



https://mlperf.org



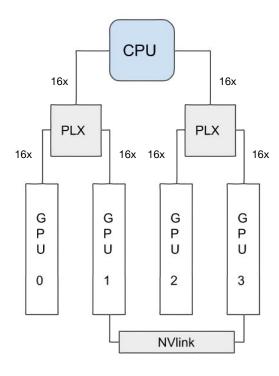
What to look for in a node / compute system

- Number of PCIe lanes. (Affects total bandwidth.)
- PCIe topology. (Tradeoff between CPU<>GPU and GPU<>GPU transfer.)
- PCIe generation. (Gen 4 is 2x bandwidth of Gen 3.)
- NUMA node topology. (Affects GPU peering & virtualization.)
- FLOPS / \$.

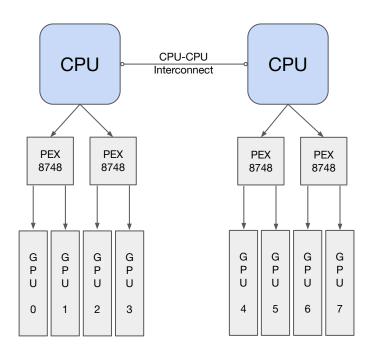
GPU peering & PCIe topology

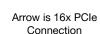


Example PCIe topology



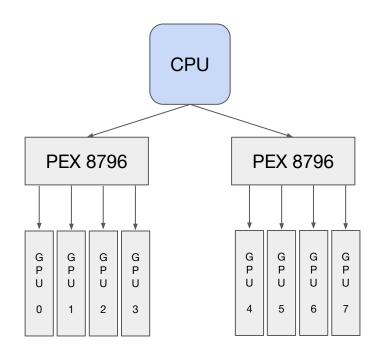
Dual root PCIe topology







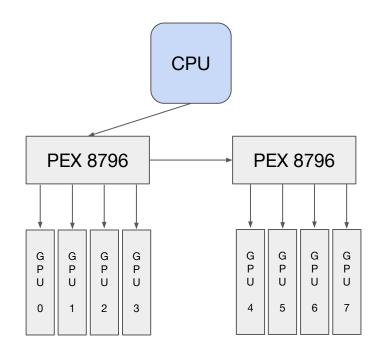
Single root PCIe topology



Arrow is 16x PCIe



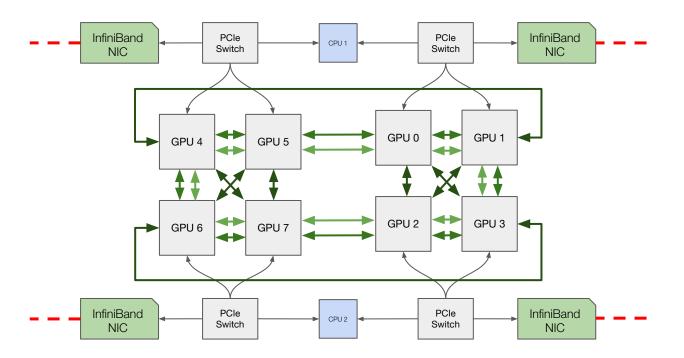
Cascaded PCIe topology

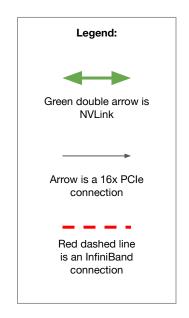


Arrow is 16x PCIe

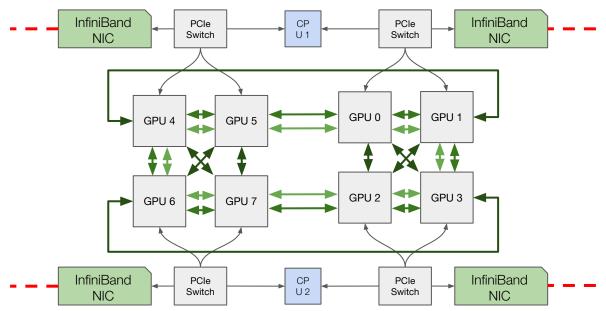


NVLink hybrid cube-mesh topology





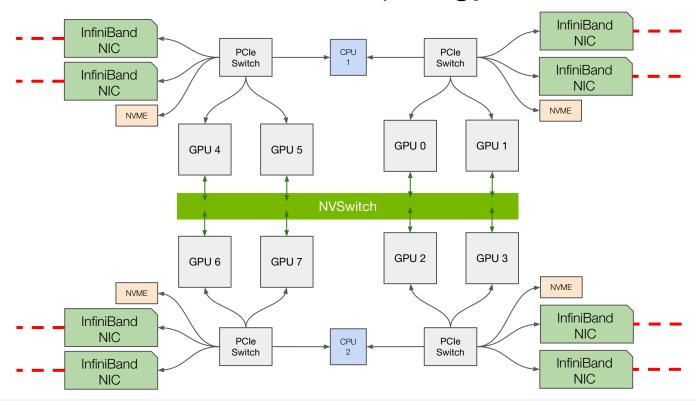
Four non-overlapping pathways



Each GPU has six NVLink connections which allows for four non-overlapping pathways to be drawn through all eight GPUs on the system and out an InfiniBand card, optimizing both GPU to GPU and node to node communication during distributed all-reduce. (See Sylvain Jeaugey's "NCCL 2.0" presentation for more information.)



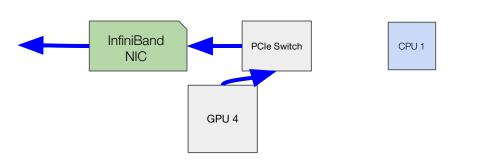
NVSwitch Topology



75



GPUDirect RDMA over InfiniBand



NIC GPU 4

PCIe Switch

InfiniBand

Data pathway with RDMA

(Directly out the door via PCIe switch)

Approximately doubles the peak node-to-node bandwidth vs. double copy.

Data pathway without RDMA

(Additional copy to CPU memory)

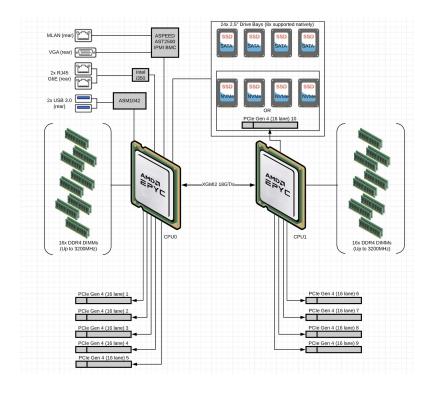


CPU 1

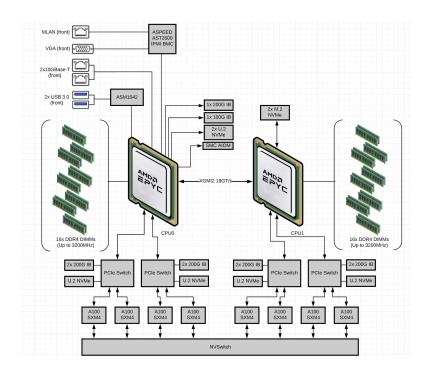
Real life examples



Lambda Scalar GPU server PCle topology



Lambda Hyperplane-8 A100 GPU server PCle topology



Putting it all together

- Write your Cluster, Rack, and Node BOMs.
- Price your BOMs.
- Negotiate a colocation contract.
- Order all of the parts on your BOMs.
- Assemble the servers.
- Ship the racks / servers to the colo / DC.
- Rack, stack, label and cable the servers.
- Install your software.
- Turn it on.

lambdalabs.com

Maintain it.

Or, work with Lambda to build a customized Lambda Echelon GPU cluster for your team.





You're done!





Citations

- Balaban, Sohmers, et al. Lambda Echelon Deep Learning GPU Cluster Reference Design Whitepaper. https://lambdalabs.com
- Sergeev, Alexander and Mike Del Balso. Horovod: fast and easy distributed deep learning in TensorFlow. https://arxiv.org/pdf/1802.05799.pdf
- Pitch Patarasuk and Xin Yuan. Bandwidth optimal all-reduce algorithms for clusters of workstations. J. Parallel Distrib. Comput., 69:117–124, 2009. https://www.cs.fsu.edu/~xyuan/paper/09jpdc.pdf
- Jeaugey, Sylvain. NCCL 2.0. (2017). http://on-demand.gputechconf.com/qtc/2017/presentation/s7155-jeaugey-nccl.pdf
- WikiChip NVLink, https://fuse.wikichip.org/news/1224/a-look-at-nvidias-nvlink-interconnect-and-the-nvswitch/
- Mellanox Cable Management Guidelines. https://www.mellanox.com/sites/default/files/products/interconnect/pdf/Mellanox Cable Management Guidelines and FAQs Application Note.pdf
- Berkeley Lab's report on ASHRAE Thermal Guidelines. https://datacenters.lbl.gov/sites/all/files/ASHRAE%20Thermal%20Guidelines %20SVLG%202015.pdf
- Mellanox CS7520 Dismantling Guide: https://www.mellanox.com/related-docs/prod ib switch systems/CS7520 Dismantling Guide.pdf
 - Additional thanks to Chuan Li, Steve Clarkson, and Thomas Sohmers for assistance with this document.



About me



- CEO of Lambda.
- Started using CNNs for face recognition in 2012.
- First employee at Perceptio. We developed image recognition CNNs that ran locally on the iPhone. Acquired by Apple in 2015.
- Research published in SPIE and NeurIPS.

About **Lambda**



Lambda provides Al infrastructure to the Fortune 500, major research institutions, and the DOD.

Our products include the Lambda Quad workstation, Lambda Hyperplane server, **Echelon GPU cluster**, and the Lambda GPU cloud.

www.lambdalabs.com

