

---

# DEEP LEARNING ON HPC SYSTEMS

**Dr. Peter Labus,**

Competence Center for High Performance Computing, Fraunhofer ITWM, Kaiserslautern.  
Fraunhofer Center Machine Learning.

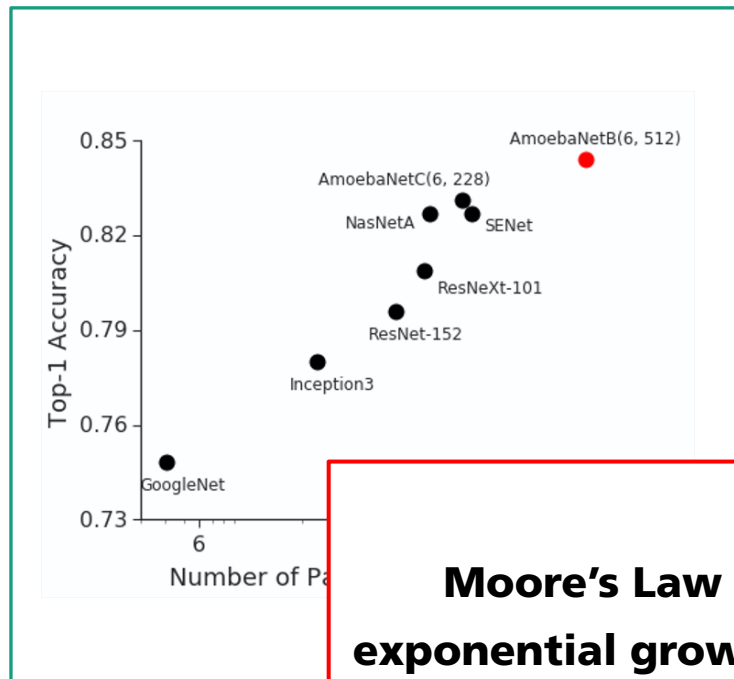
---



**9<sup>th</sup> HPC Status conference  
Gauß-Allianz  
October 17, 2019,  
Paderborn**

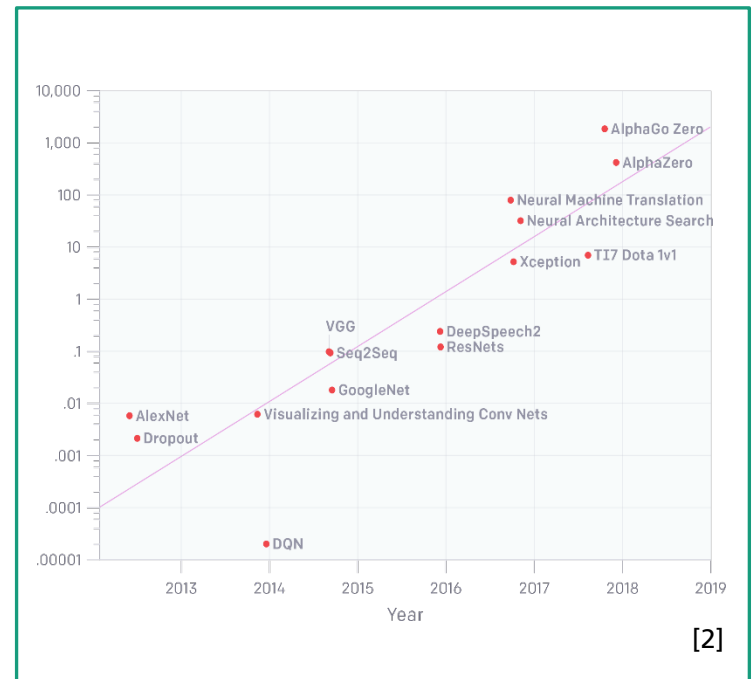
# Increasing AI model complexity leads to increasing compute demand

## increasing AI model complexity



**Moore's Law of AI:  
exponential growth with a  
3.5 month-doubling time**

## increasing compute demand



[1] Huang, Yanping, et al. "Gpipe: Efficient training of giant neural networks using pipeline parallelism."

[2] <https://openai.com/blog/ai-and-compute/>

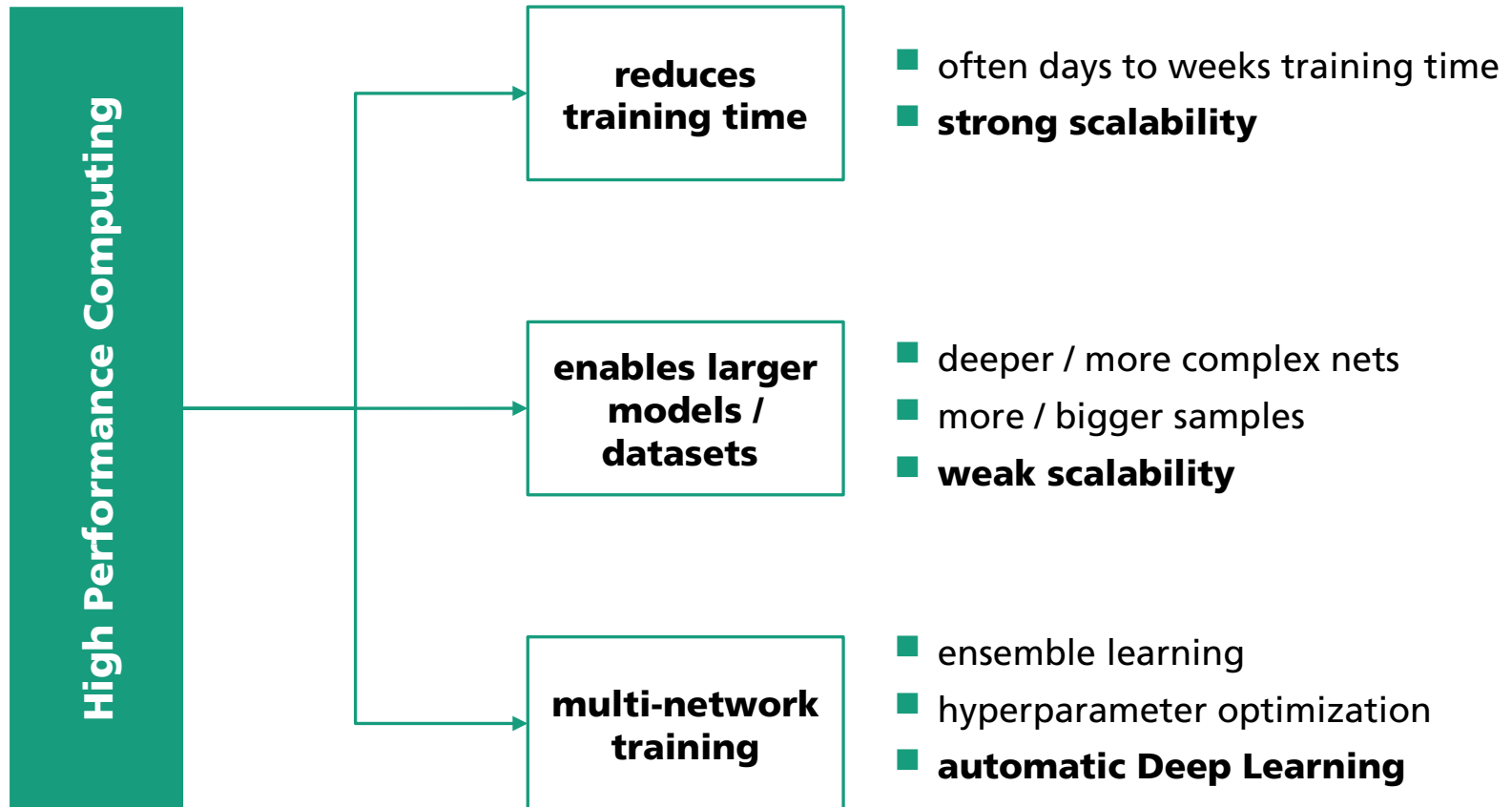
---

# AGENDA

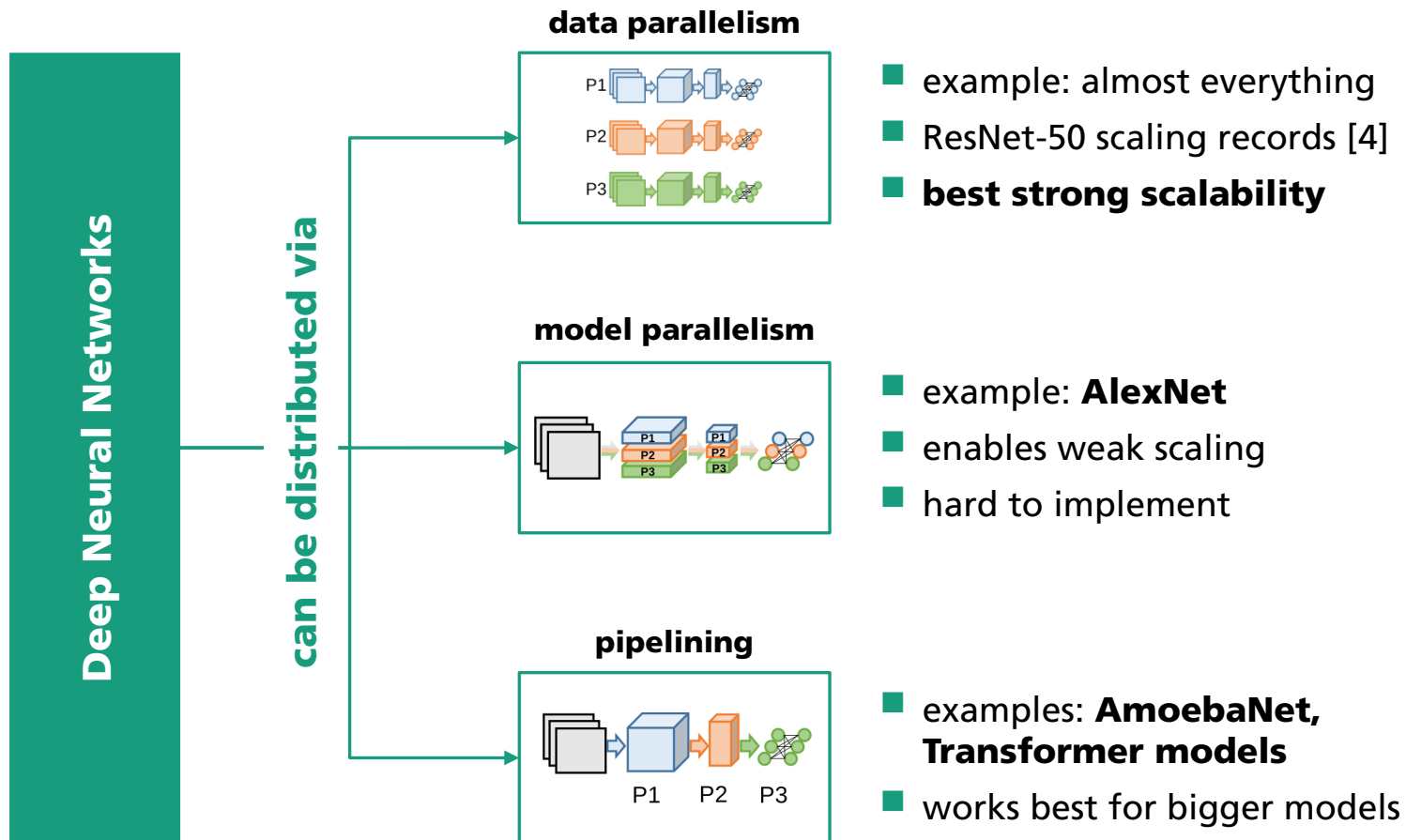
---

1. Deep Learning & High Performance Computing
2. Scaling Deep Neural Network Training
3. Novel algorithms & visualization
4. Automatic Deep Learning

# HPC reduces training time, enables larger models and datasets & multi-network training



# DNNs can be distributed via data parallelism, model parallelism & pipelining

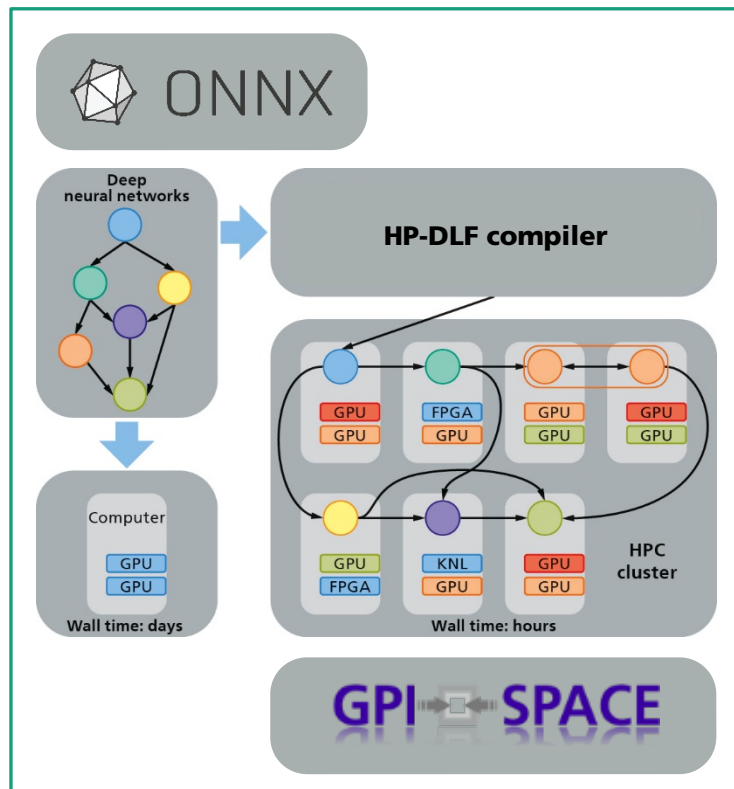


[3] Ben-Nun, Tal, and Torsten Hoefer. "Demystifying parallel and distributed deep learning" ACM Computing Surveys (CSUR) 52.4 (2019):65.

[4] Mikami, Hiroaki, et al. "Imagenet/resnet-50 training in 224 seconds." arXiv preprint arXiv:1811.05233 (2018). C

# Initial design of HP-DLF led to a number of issues

## initial design of HP-DLF



## a number of issues

### ■ ONNX as frontend

no training interface (solvers, loss functions, data processing), not widely supported, user needs to learn new interface

virtualization overhead, need more fine-grained control

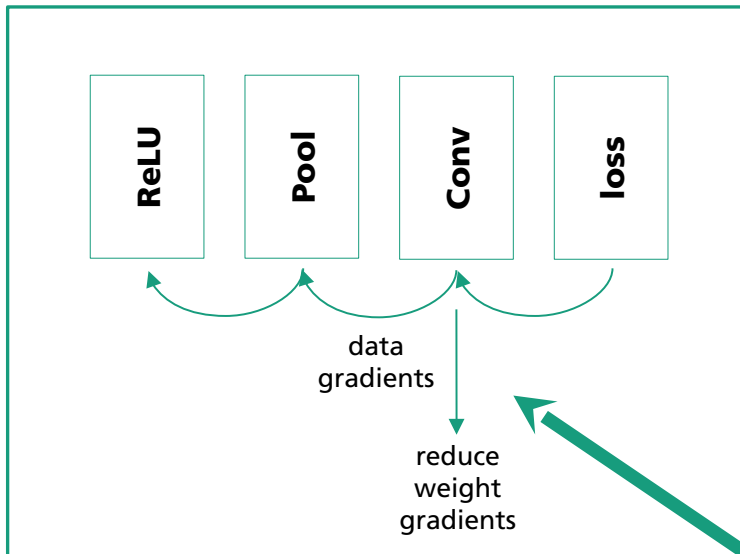
**Build scalable communication library on top of TensorFlow**

### ■ framework development

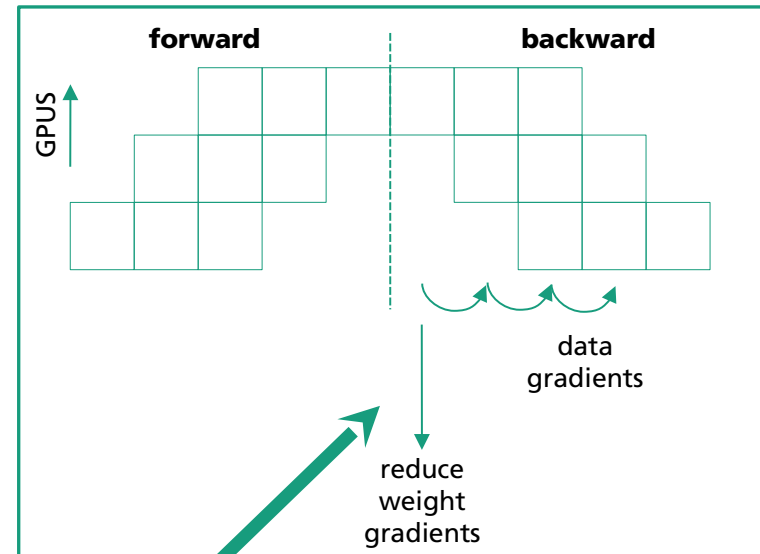
lots of development overhead for user interface, data processing, implementation wrappers, ...

# GASPI communication library enables the overlap of computation & communication

## data-parallel SGD



## pipelined SGD



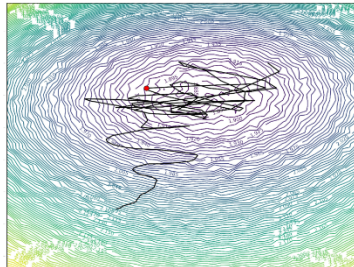
**enables**

**overlap computation & communication**

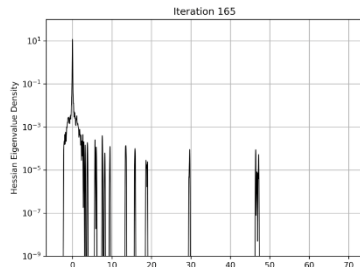
[5] Huang, Yanping, et al. "Gpipe: Efficient training of giant neural networks using pipeline parallelism." arXiv preprint arXiv:1811.06965 (2018).

# New visualization techniques help study novel optimizers

## New visualization techniques



2D/3D loss surfaces & trajectories



Hessian eigenvalue spectra

### Our contributions: [6]

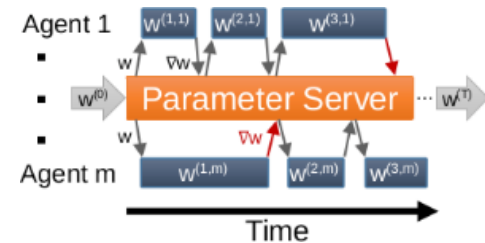
- projection onto EVs of Hessian
- scalable visualization impl.
- scalable & efficient Lanczos

[github.com/cc-hpc-itwm/GradVis](https://github.com/cc-hpc-itwm/GradVis)

help  
to  
study

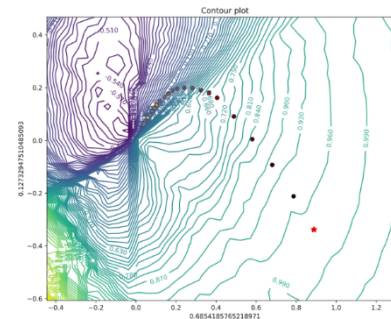
## Novel optimizers

### Asynchronous SGD [7]



improve scaling by reduction of message frequency

### ResOpt [8]



improve scaling by reduction of message sizes

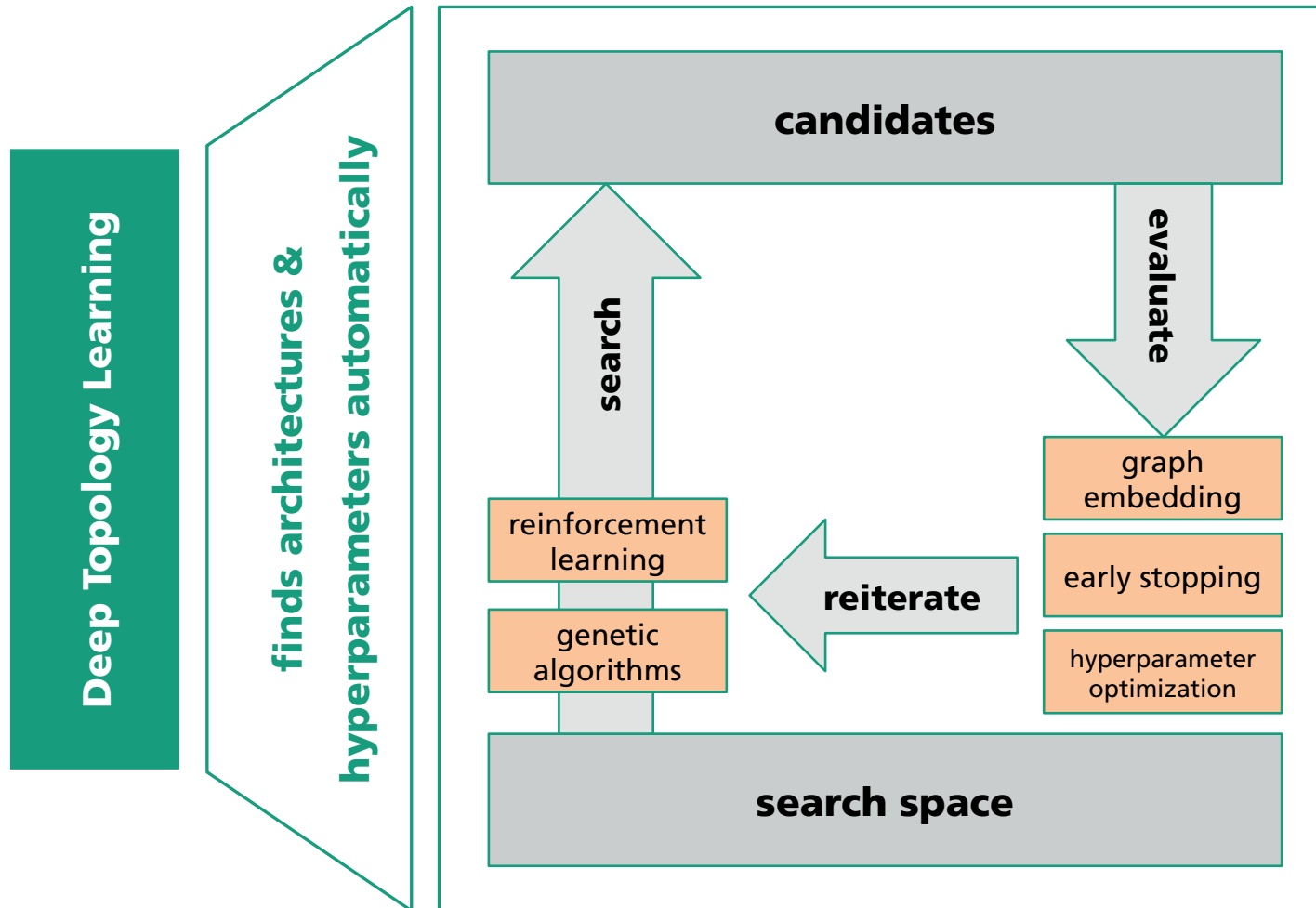
[6] Chatzimichailidis, A., et al. "GradVis: Visualization and Second Order [...]" arXiv preprint arXiv:1909.12108 (2019).

[7] Keuper, J. et al. "Asynchronous parallel stochastic gradient descent [...]" Proceedings of the Workshop on Machine Learning in High-Performance Computing Environments. ACM, 2015.

[8] Lorch, D. "ResOpt [...]" to appear

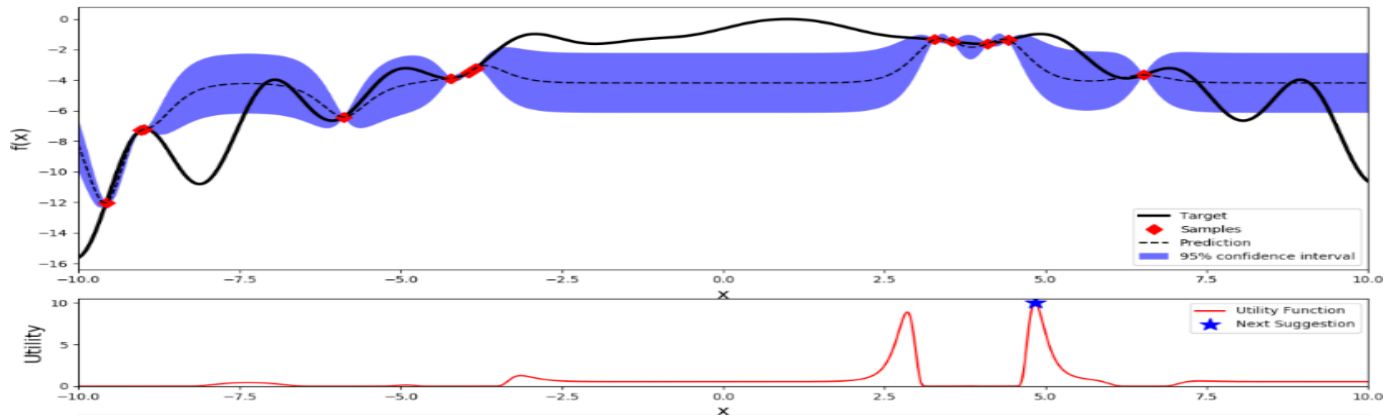


# Deep Topology Learning finds DNN architectures & hyperparameters automatically



# Lazy Gaussian Processes enable scalable Bayesian Optimization

Lazy Gaussian Processes



Build covariance matrix iteratively: 
$$K_{n+1} = \begin{pmatrix} K_n & \mathbf{p} \\ \mathbf{p}^T & c \end{pmatrix}$$

enables

Scalable Bayesian Opt

**Matrix inversion via Cholesky decomposition: improvement from  $O(n^3)$  to  $O(n^2)$**

Naïve Cholesky decomposition	
Iteration	Accuracy
16	0.74
26	0.75
43	0.77
50	0.78
176	0.79

Optimized Cholesky decomposition	
Iteration	Accuracy
1	0.11
16	0.51
21	0.77
35	0.79
61	<b>0.80</b>

[6] Ram, Raju, et. al. "Scalable Hyperparameter Optimization with Lazy Gaussian Processes" to appear in Proceedings of the Workshop on Machine Learning in High-Performance Computing Environments. ACM, 2019.

# Summary: Deep Learning on HPC systems

## Insight

