



Sussex, Monday 3 July 2023

Can we learn from the brain to make AI more energy efficient?

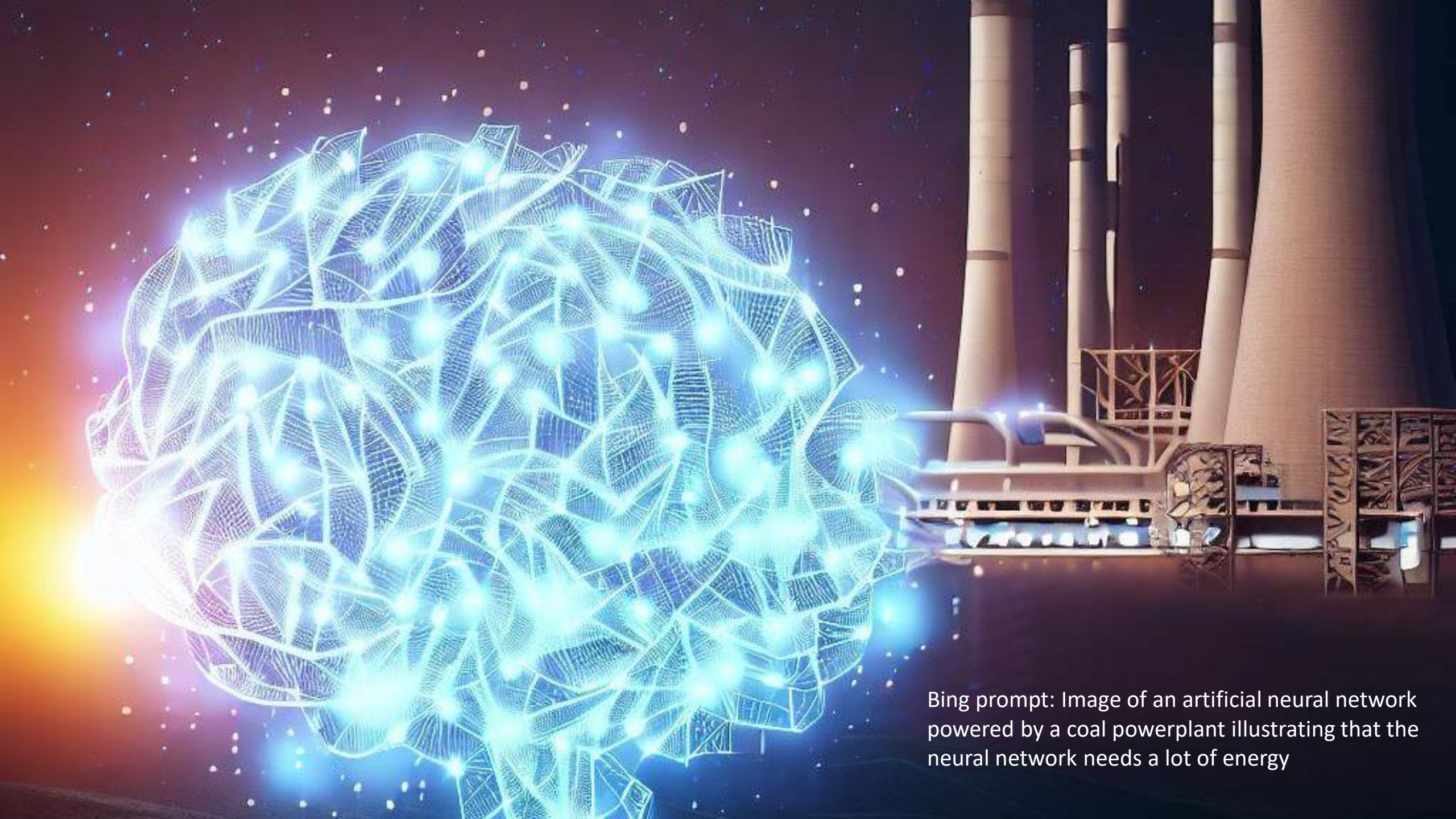
Prof Thomas Nowotny

School of Engineering and Informatics

University of Sussex

Bing prompt: A Spiking Neural Network



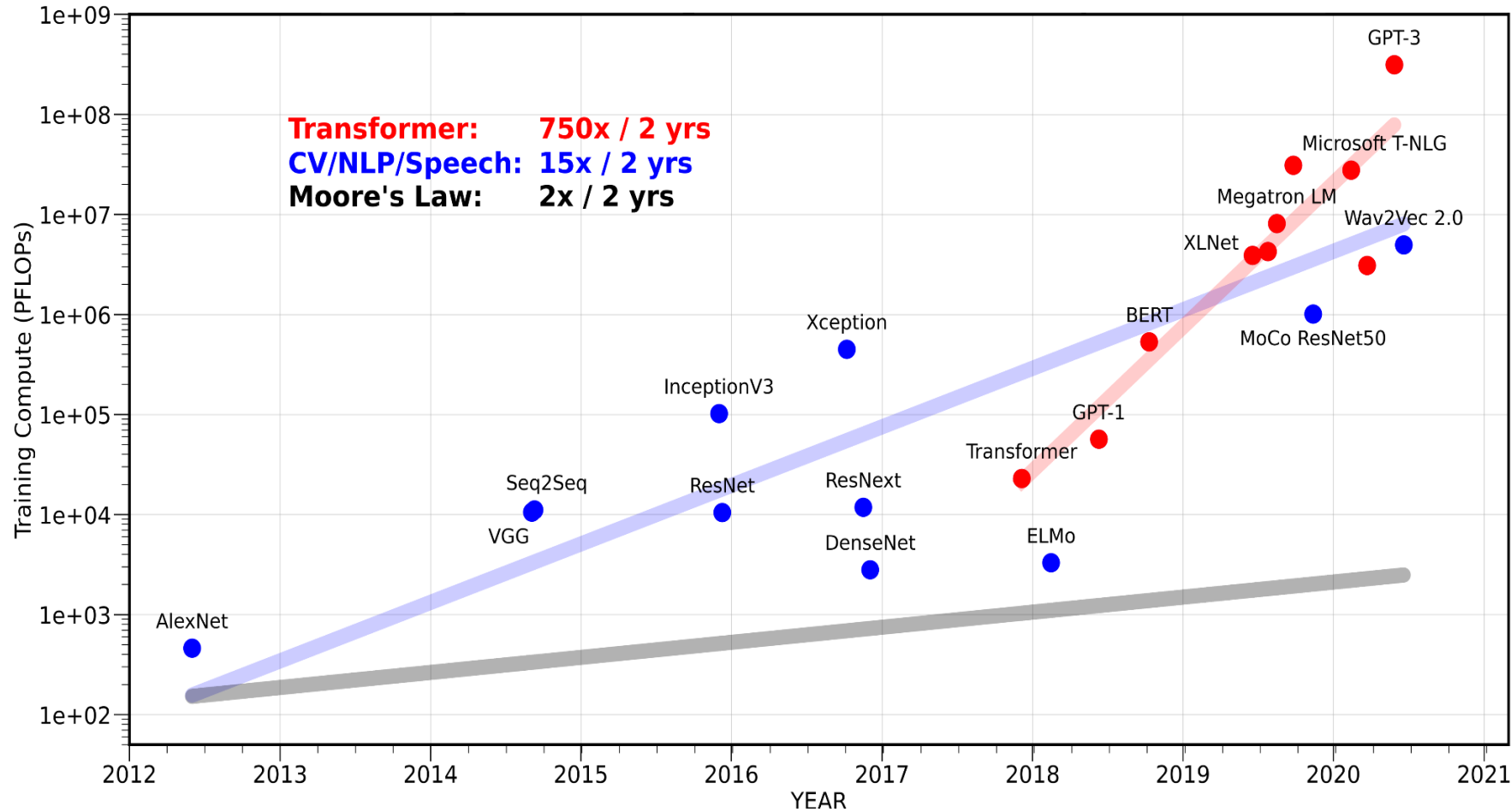


Bing prompt: Image of an artificial neural network powered by a coal powerplant illustrating that the neural network needs a lot of energy



Bing prompt: Car driving to the moon, using a lot of petrol

Modern AI Resource Trends



Gholami, A., Yao, Z., Kim, S., Mahoney, M. W., & Keutzer, K. (2021) AI and Memory Wall. *RiseLab Medium* Blog Post, University of California Berkeley, 2021, March 29





Brain Circuits vs Artificial Neural Networks

- Electro-chemical “wetware”
- Evolved embodied
- Sparse communication with “spikes”
- Metal-oxide hardware
- Built for machine learning
- Dense communication of floating point numbers

Neuromorphic computing



SpiNNaker
(Manchester/
Dresden)

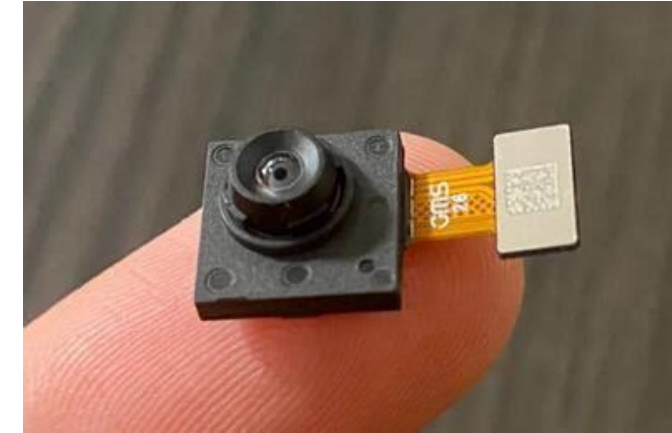


BrainScaleS
(Heidelberg)



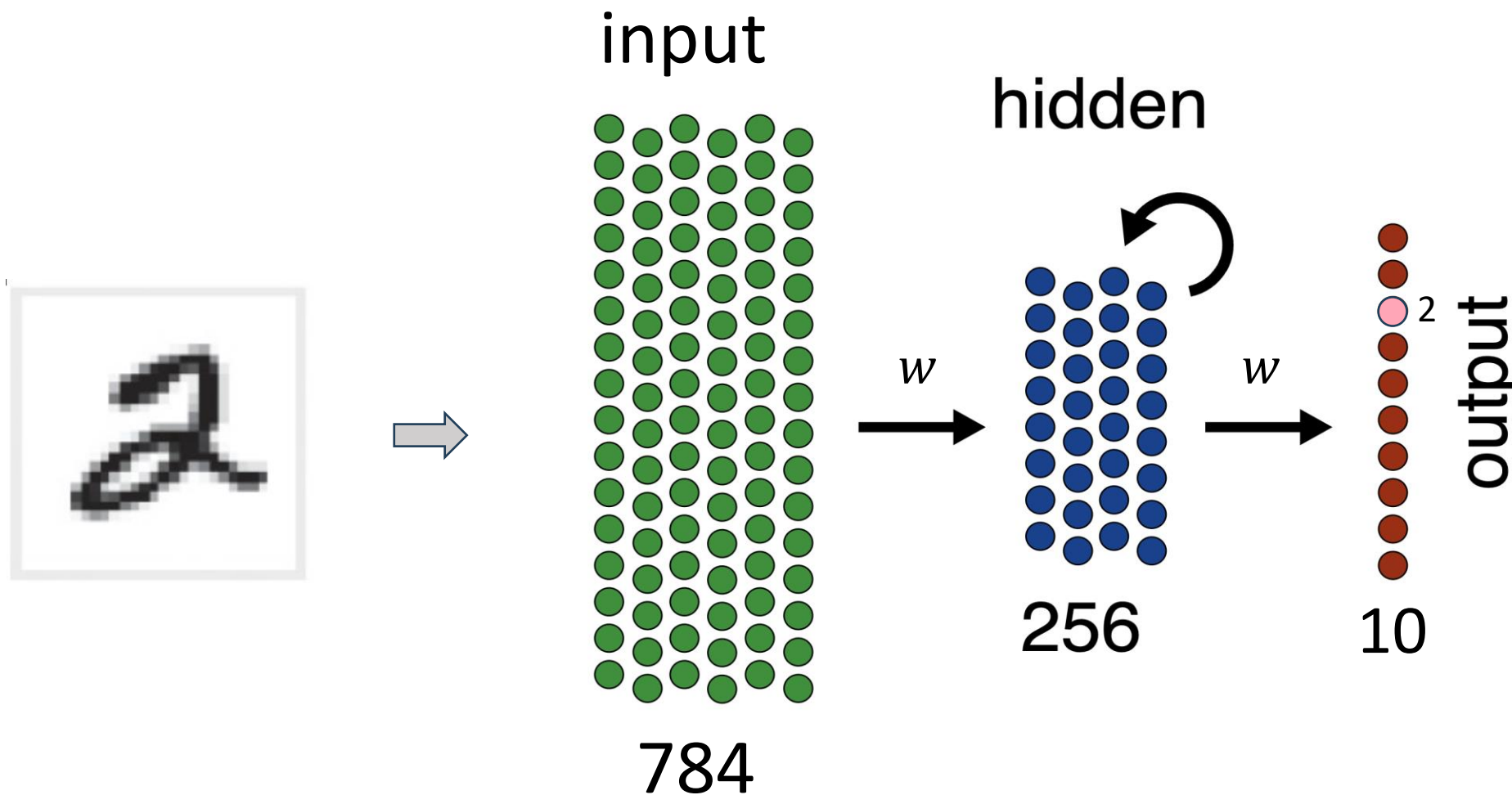
Intel Loihi
(Intel Corporation)

Synsense
Speck

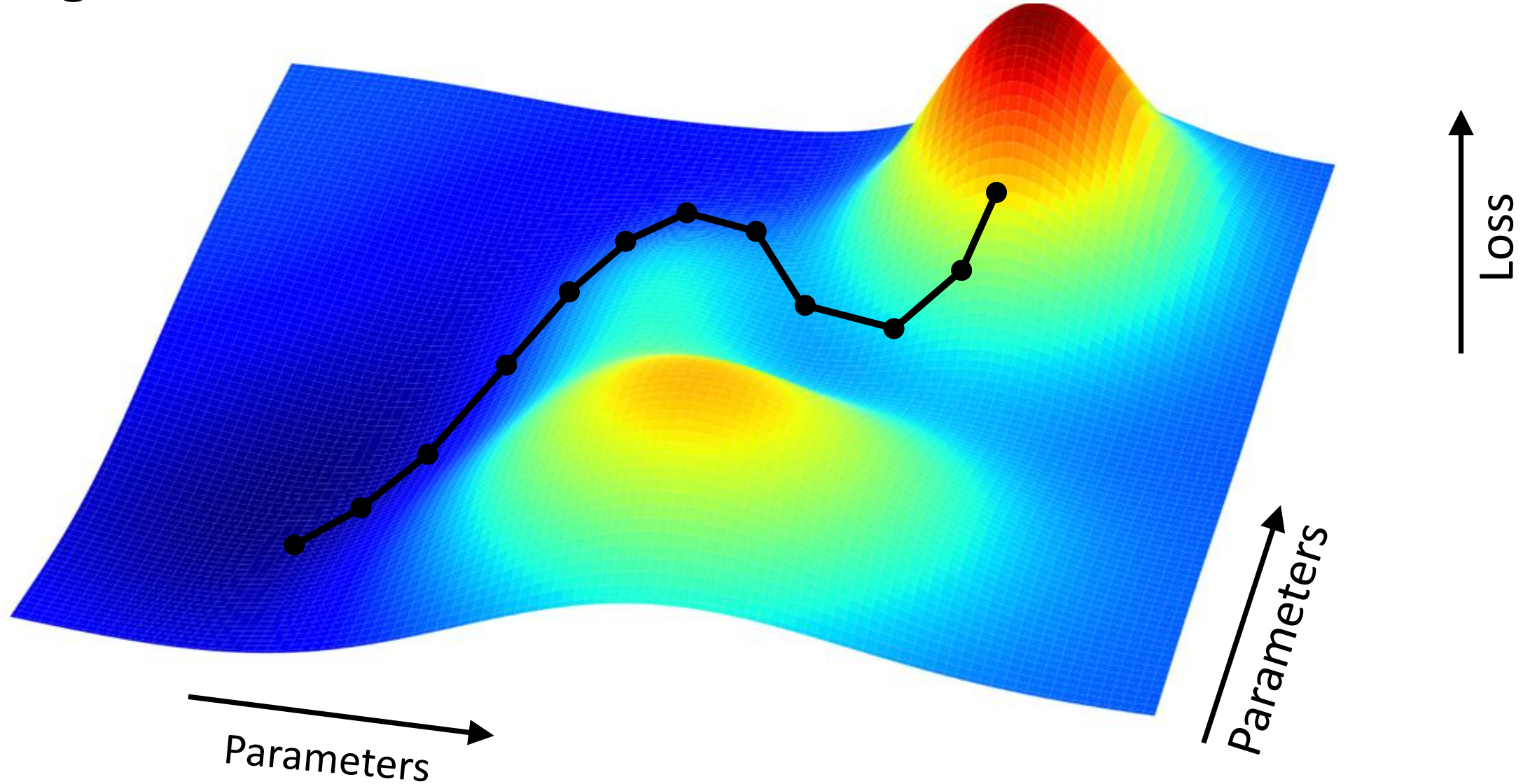


IBM synapse
(IBM)

General classification problem



Training: Gradient descent



IOP Publishing

Neuromorph. Comput. Eng. 2 (2022) 024002

<https://doi.org/10.1088/2634-4386/ac5ac5>

NEUROMORPHIC

Computing and Engineering

PAPER



mlGeNN: accelerating SNN inference using GPU-enabled neural networks

OPEN ACCESS

James Paul Turner^{1,*} , James C Knight¹ , Ajay Subramanian² and Thomas Nowotny¹ 

RECEIVED
22 December 2021

REVISED
9 February 2022

ACCEPTED FOR PUBLICATION
4 March 2022

PUBLISHED
25 March 2022

¹ Centre for Computational Neuroscience and Robotics, School of Engineering and Informatics, University of Sussex, Brighton, United Kingdom

² Department of Psychology, New York University, New York, NY 10003, United States of America

* Author to whom any correspondence should be addressed.

E-mail: J.P.Turner@sussex.ac.uk

Keywords: machine learning, spiking neural networks, GPU, ANN to SNN conversion, convolutional neural networks, GeNN, Re

Original content from this work may be used under the terms of the

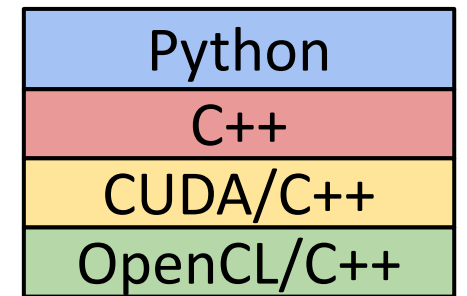
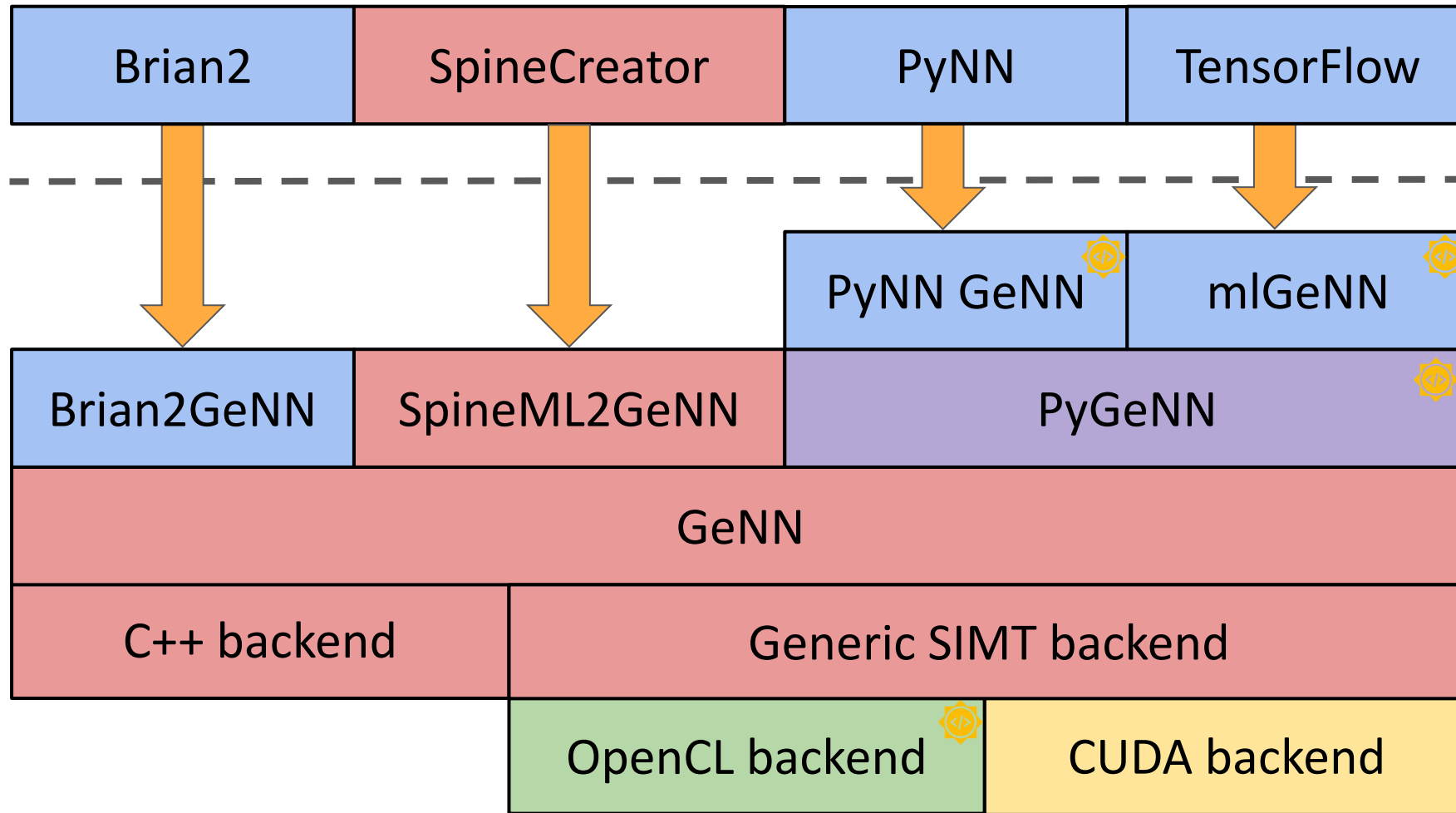


GeNN (**G**PU enhanced **N**eural **N**etworks)

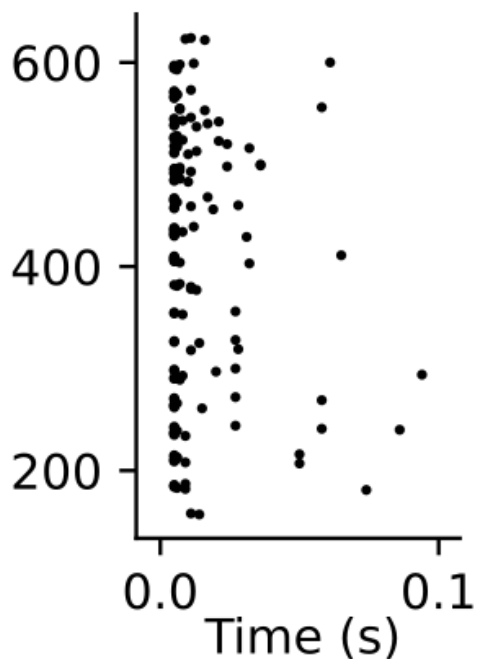
GeNN generates CUDA kernels and data transfer “convenience functions” based on a model definition provided by a user.

(Meta Compiler)

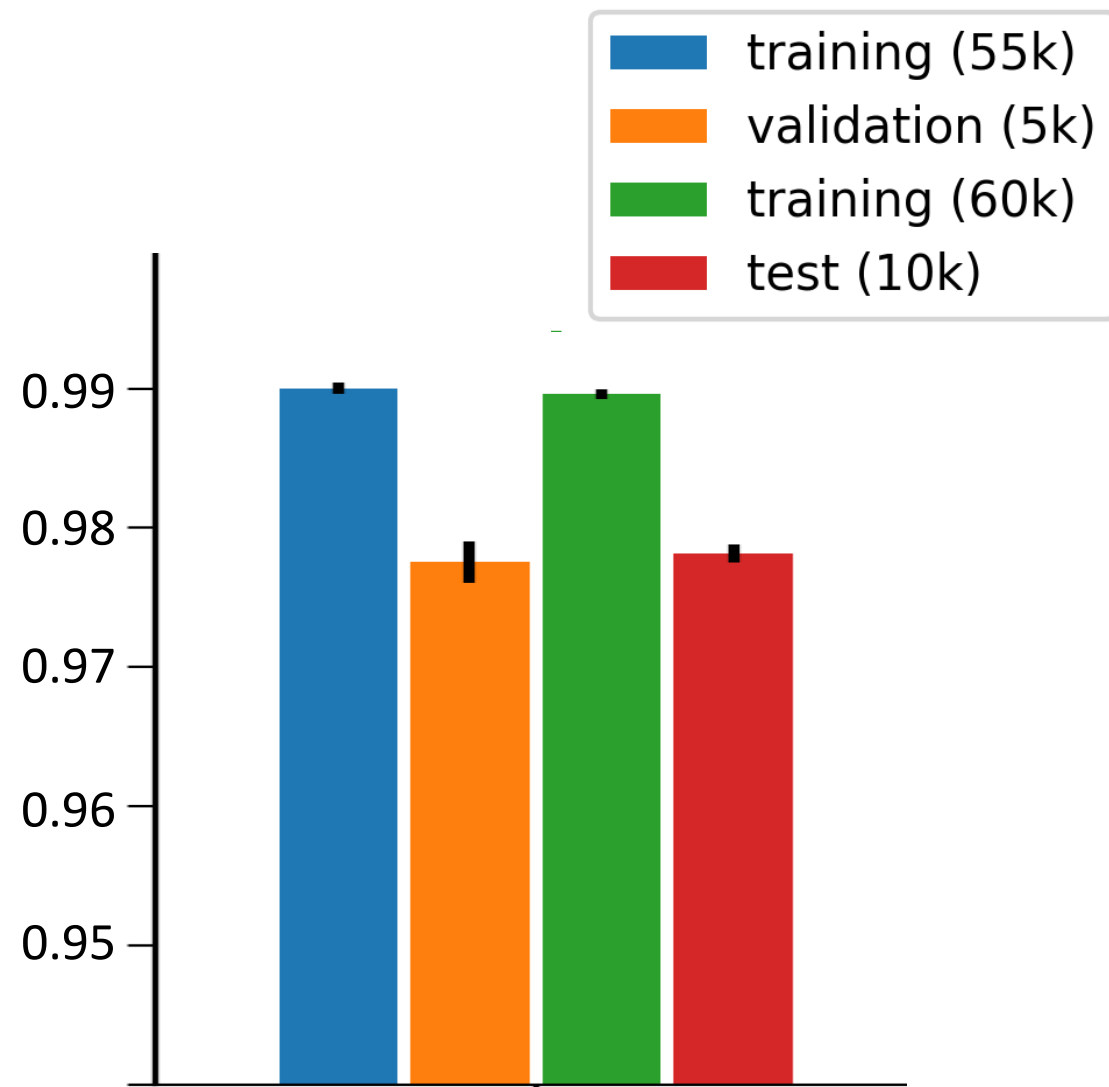
GeNN ecosystem



Handwritten digits with EventProp



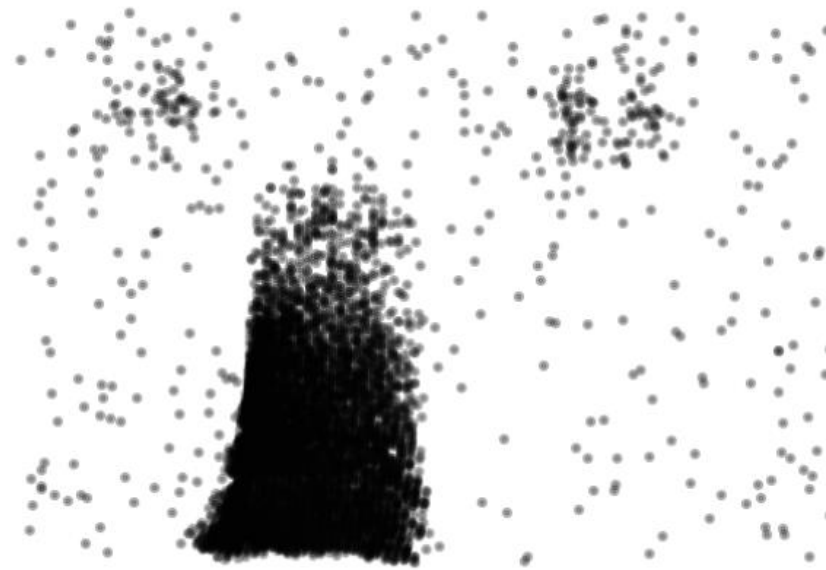
Trains to >97% in 27s on a PC with RTX 3090



Spiking Heidelberg Digits



"three"

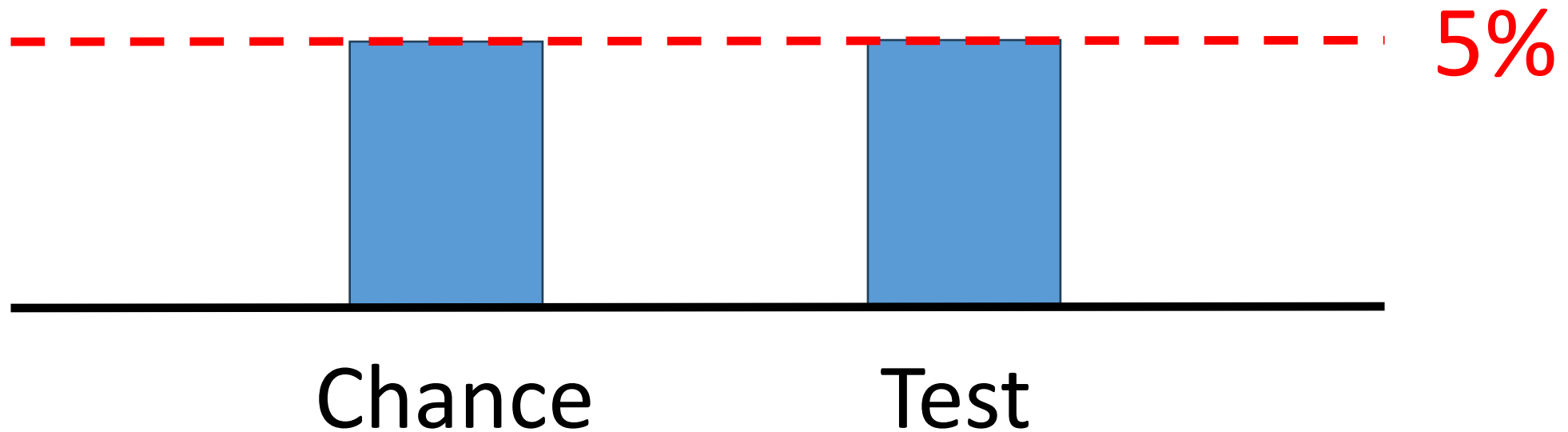


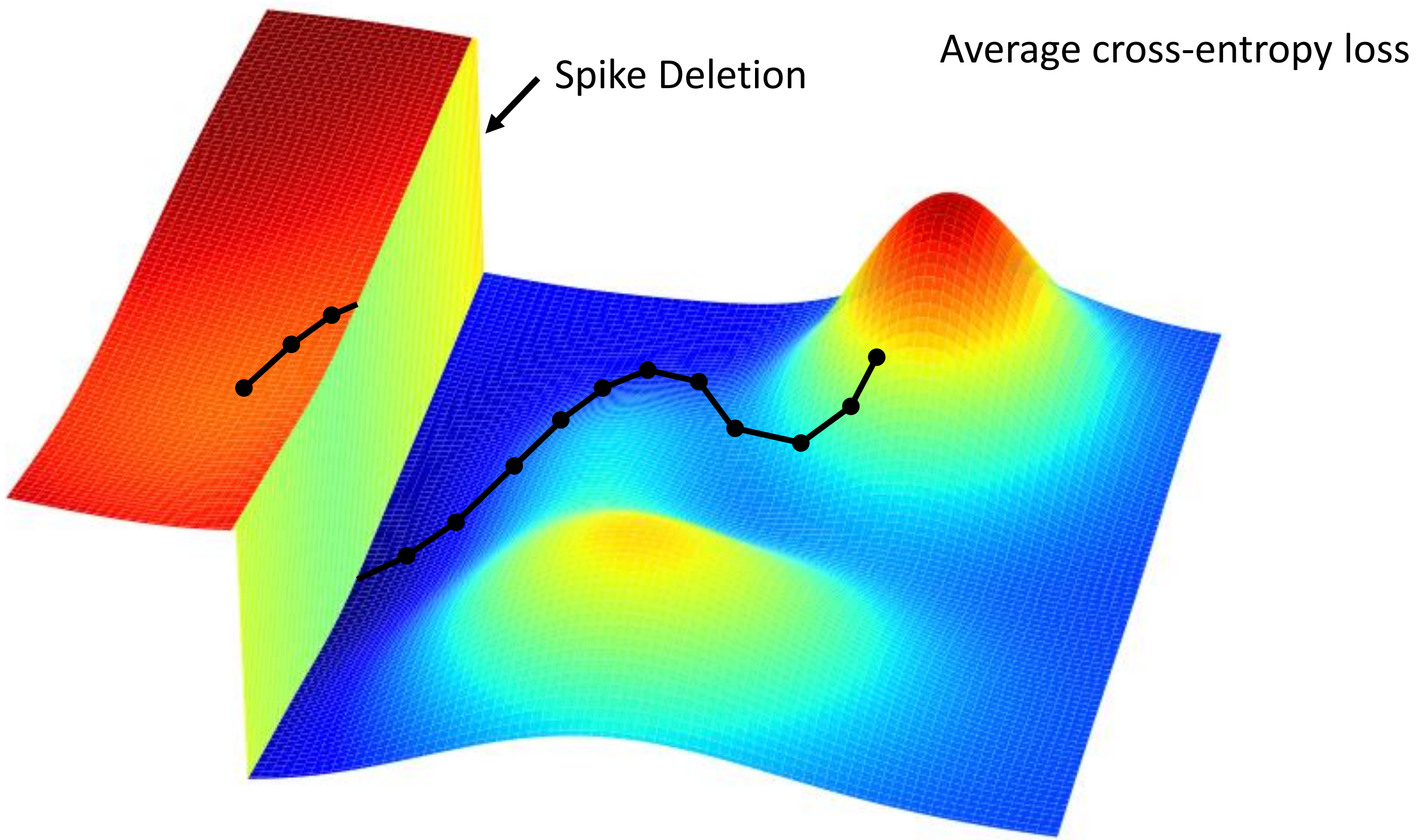
"seven"

Cramer, B., Stradmann, Y., Schemmel, J., & Zenke, F. (2020). The Heidelberg Spiking Data Sets for the Systematic Evaluation of Spiking Neural Networks. *IEEE Transactions on Neural Networks and Learning Systems*.

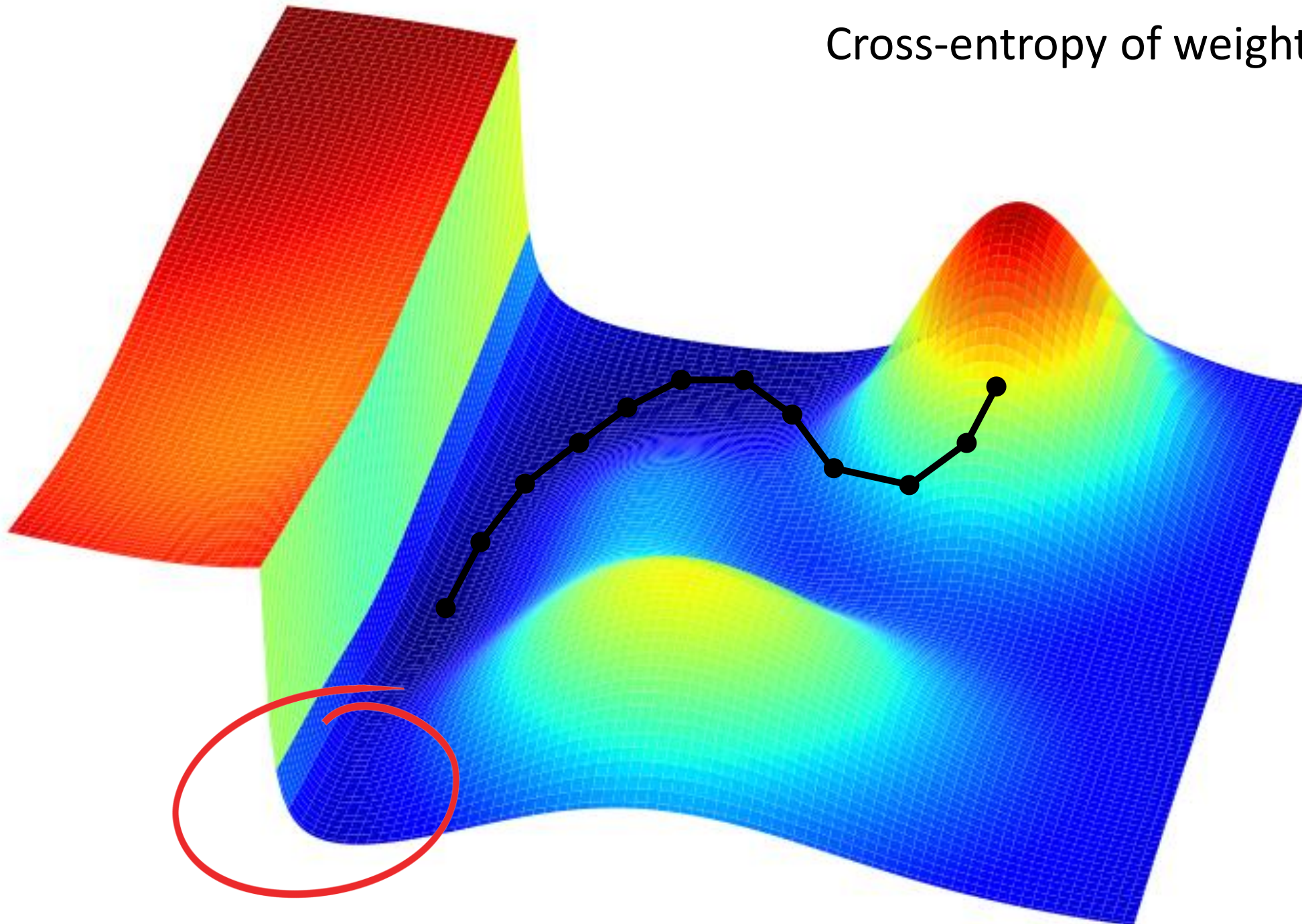
<https://doi.org/10.1109/TNNLS.2020.3044364>

No learning!

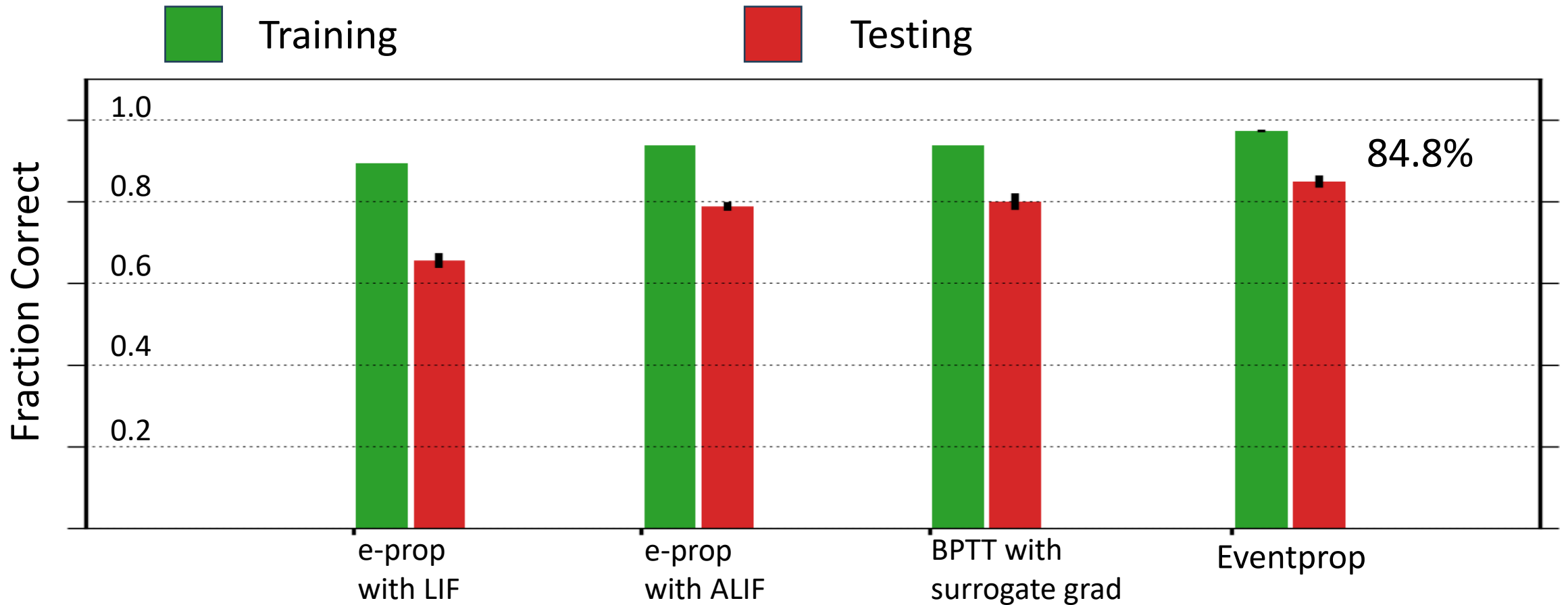




Cross-entropy of weighted average loss



Spiking Heidelberg Digits



1. Zenke F, Vogels T. The remarkable robustness of surrogate gradient learning for instilling complex function in spiking neural networks. *Neural Comput* 33:899-925 (2021)
2. J. C Knight, T. Nowotny. Efficient GPU training of LSNNs using eProp. *NICE 2022. ACM*, 8–10



EventProp
Spiking simulation

Gradient Descent

Augmentation

Augmentation

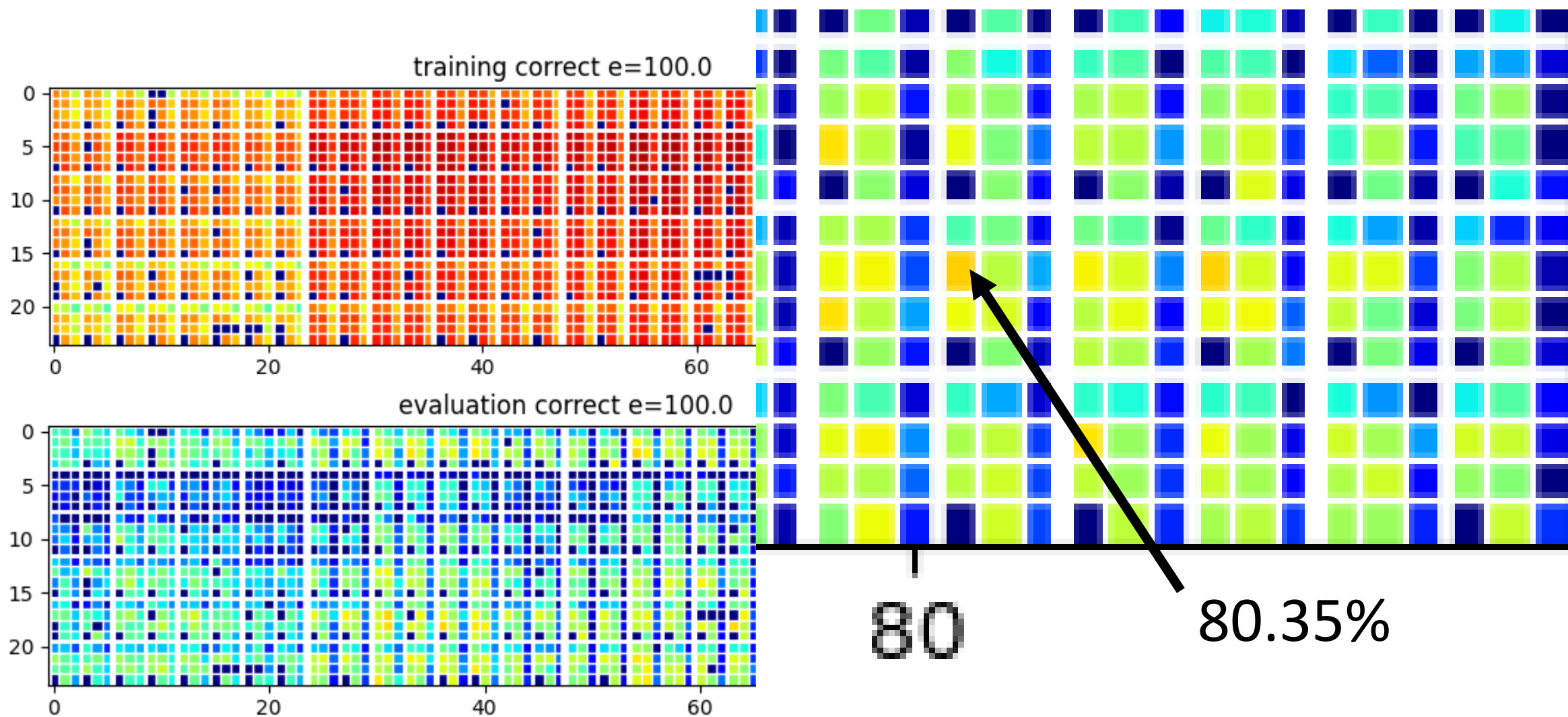
- Augmentation is used to make data sets larger to counter-act overfitting to few examples
- Augmentation for event-based data is an open field with not much work
- Candidates:
 - Global random channel shift
 - Local random channel swap
 - Random time compression/dilation
 - Blending

Parameter sweeps

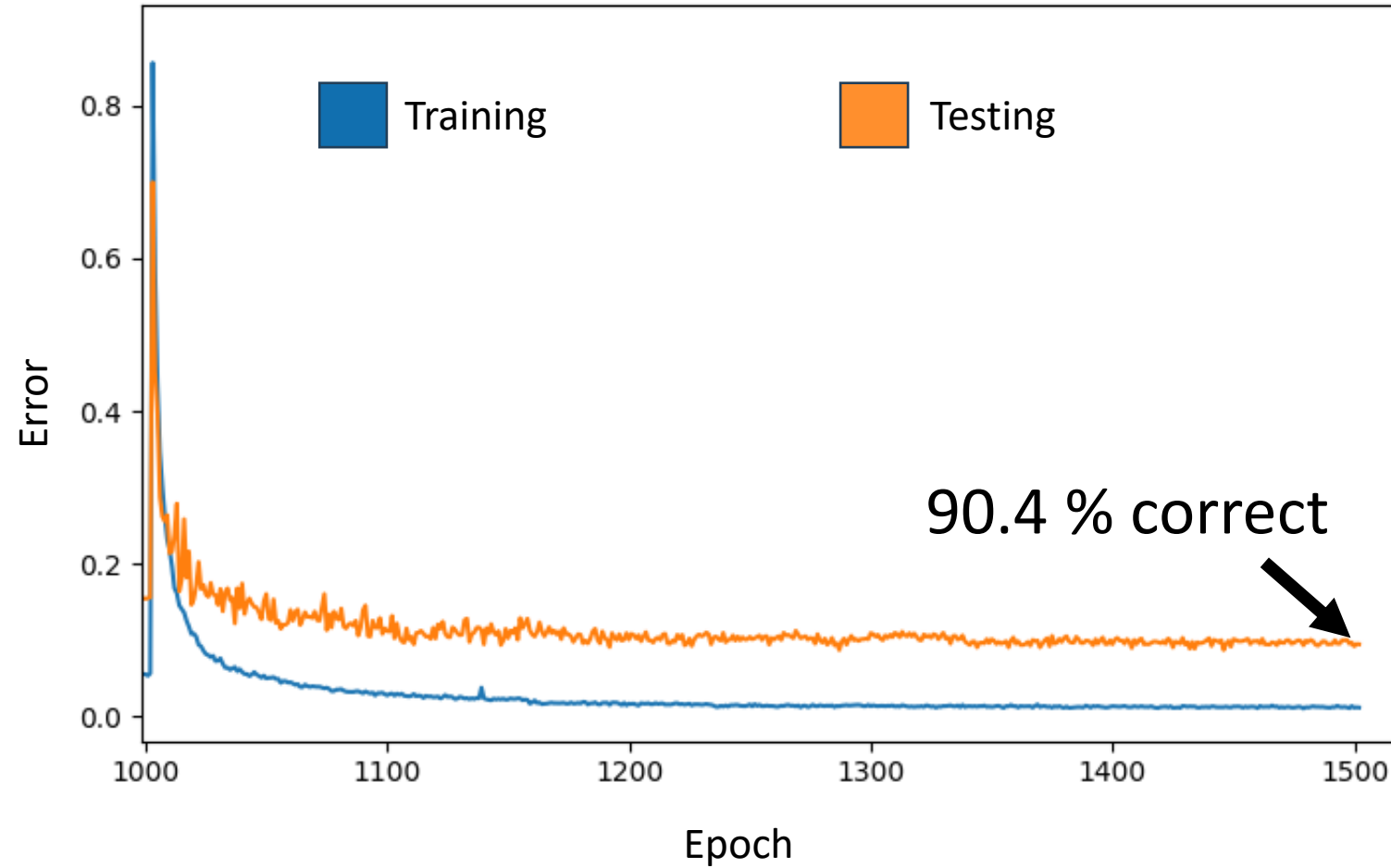
- Try parameters for the network and augmentations
- Best proxy for test set performance is leave-one-speaker-out (loso) cross-validation
- A full 10-fold loso cross-validation with 300 epochs per fold takes about **16 hours on an A100 GPU**



Example Scan Results



Test set performance



Summary

- JADE 2 (and other GPU clusters) help us investigate event-based neural network solutions competitive against SOTA

Dataset	Method	Recurrent	Delays	# Params	Top1 Accuracy
SHD	EventProp-GeNN [31]	✓	✗	N/a	84.80 ± 1.5%
	Cuba-LIF [7]	✗	✗	N/a	87.80 ± 1.1%
	Adaptive SRNN [44]	✓	✗	N/a	90.40% ← 90.4 % correct
	SNN with Delays [2]	✗	✓	0.1M	90.43%
	TA-SNN [43]	✗	✗	N/a	91.08%
	STSC-SNN [46]	✗	✗	2.1M	92.36%
	Adaptive Delays [37]	✗	✓	0.1M	92.45%
	RadLIF [3]	✓	✗	3.9M	94.62%
Our work (2 hidden layers)	✗	✓	0.2M	95.07 ± 0.24% } Use test set for validation	

Hammouamri, I., Khalfaoui-Hassani, I., & Masquelier, T. (2023). Learning Delays in Spiking Neural Networks using Dilated Convolutions with Learnable Spacings. *arXiv preprint arXiv:2306.17670*.

Acknowledgements

- Jamie Knight
- Other GeNN developers:
 - Esin Yavuz
 - James Turner
 - Anton Komissarov
 - Fawad Ali, Obaid-Ur-Rehman, Muhammad Irfan Ali, University of Engineering and Technology, Taxila, PAKISTAN (OpenCL backend - internship)

