



Hardware architecture and routing-aware training for optimal memory usage: a case study



A work of

Theo Ballet and Jimmy Weber

Supervised by

Melika Payvand





ETH zürich

institute of neuroinformatics

















Hardware-software co-design









Hardware-software co-design $\min_{\theta, A} \mathbb{E}[\mathcal{L}(\theta) | C(\theta, A, R_A) \leq 0]$

...But first, let's get inspired from another method...

Problem:

"Train a neural network with fixed sparsity (e.g. 7 weights) and unconstrained connectivity."

Solution:

Dynamical architecture search





 $|\theta|_0$: Number of non-zero elements 5 $|\theta|_1$: Sum of all elements

1. BELLEC, Guillaume, KAPPEL, David, MAASS, Wolfgang, et al.Deep rewiring: Training very sparse deep networks. arXiv preprint arXiv:1711.05136, 2017.







Why is Jimmy telling you about this?





Hardware-software co-design $\min_{\theta, A} \mathbb{E}[\mathcal{L}(\theta) | C(\theta, A, R_A) \leq 0]$

...But first, let's get inspired from another method...

Problem:

"Train a neural network with fixed sparsity (e.g. 7 weights) and unconstrained connectivity."

Solution:

Dynamical architecture search



DeepR¹





 $|\theta|_0$: Number of non-zero elements 7 $|\theta|_1$: Sum of all elements

1. BELLEC,G, et al. Deep rewiring: Training very sparse deep networks. 2017.





 $|\theta|_0$: Number of non-zero elements 8 $|\theta|_1$: Sum of all elements





 $|\theta|_0$: Number of non-zero elements 9 $|\theta|_1$: Sum of all elements





Hardware-software co-design $\min \mathbb{E}[\mathcal{L}(\theta) | C(\theta, A, R_A) \leq 0]$ θ . A

Extension of DeepR



> Evaluating $C(\theta)$ should be efficient.

```
▶ Get \gamma(\theta), a proxy of C(\theta).
```



The Mosaic² as a case study



Neuron Tile (0,0)



The Mosaic² as a case study





2: Dalgaty, T. et al. Mosaic: in-memory computing and routing for small-world spike-based neuromorphic systems. 2024



The Mosaic² as a case study





Constraints of MOSAIC mappability = routing resources























Sparsity level of inter-core distance $p_{d=3}$

 $\widehat{P} = \{\widehat{p_1}, \widehat{p_3}, \widehat{p_5}, \dots\}$

Inter-core sparsity level, for core at distance d.









With: $|\theta|_0$: Number of non-zero elements $|\theta|_1$: Sum of all elements

19



 $\gamma(\theta) = |\theta_d|_1$

$$\widehat{P} = \{\widehat{p_1}, \widehat{p_3}, \widehat{p_5}, \dots\}$$

Inter-core sparsity level, for core at distance d.







Hardware-software co-design













Hardware-software co-design $\min_{\theta, A} \mathbb{E}[\mathcal{L}(\theta) | C(\theta, A, R_A) \leq 0]$



Further insights and outlook



ightarrow A family of HW have the same accuracy

→ Accuracy is a function of sparsity profile (inter-population sparsity based on their distance)



Take home message



Mathematical formulation of architectural and routing constrains: $C(\theta, A, R_A)$ Developing a method inspired by DEEPR $\min_{\theta} \mathbb{E}[\mathcal{L}(\theta) | C(\theta, A, R_A) \leq 0]$

Hardware architecture and routing-aware training for optimal memory usage: a case study

Such that C = 0

The Mosaic Architecture

⇒ Get accuracy/memory improvement

Acknowledgment



The EIS-lab, retreat 2024



Théo Ballet, Paris.



Melika Payvand, ETH/UZH