

OnDevice Learning

Hahn-Schickard
Marcus.rueb@hahn-schickard.de

What is OnDevice Training and why is it needed?

Advantages [1]:



Privacy



Continuously Learning



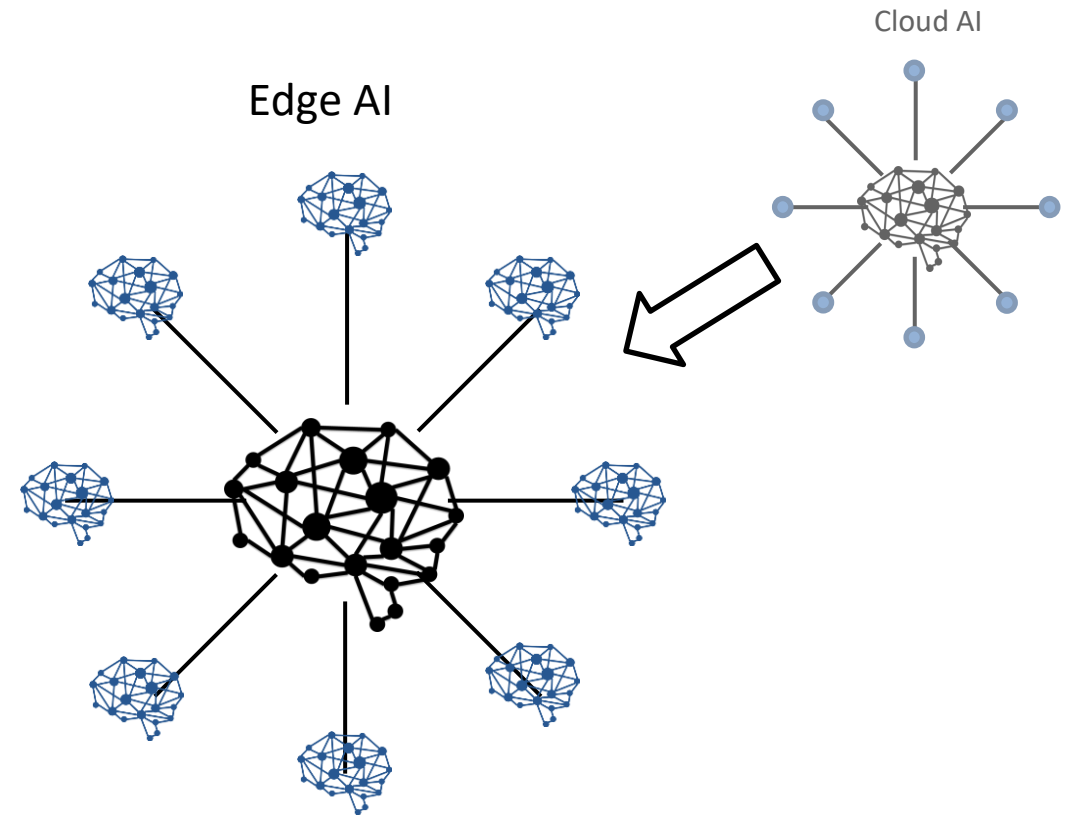
Energy-saving



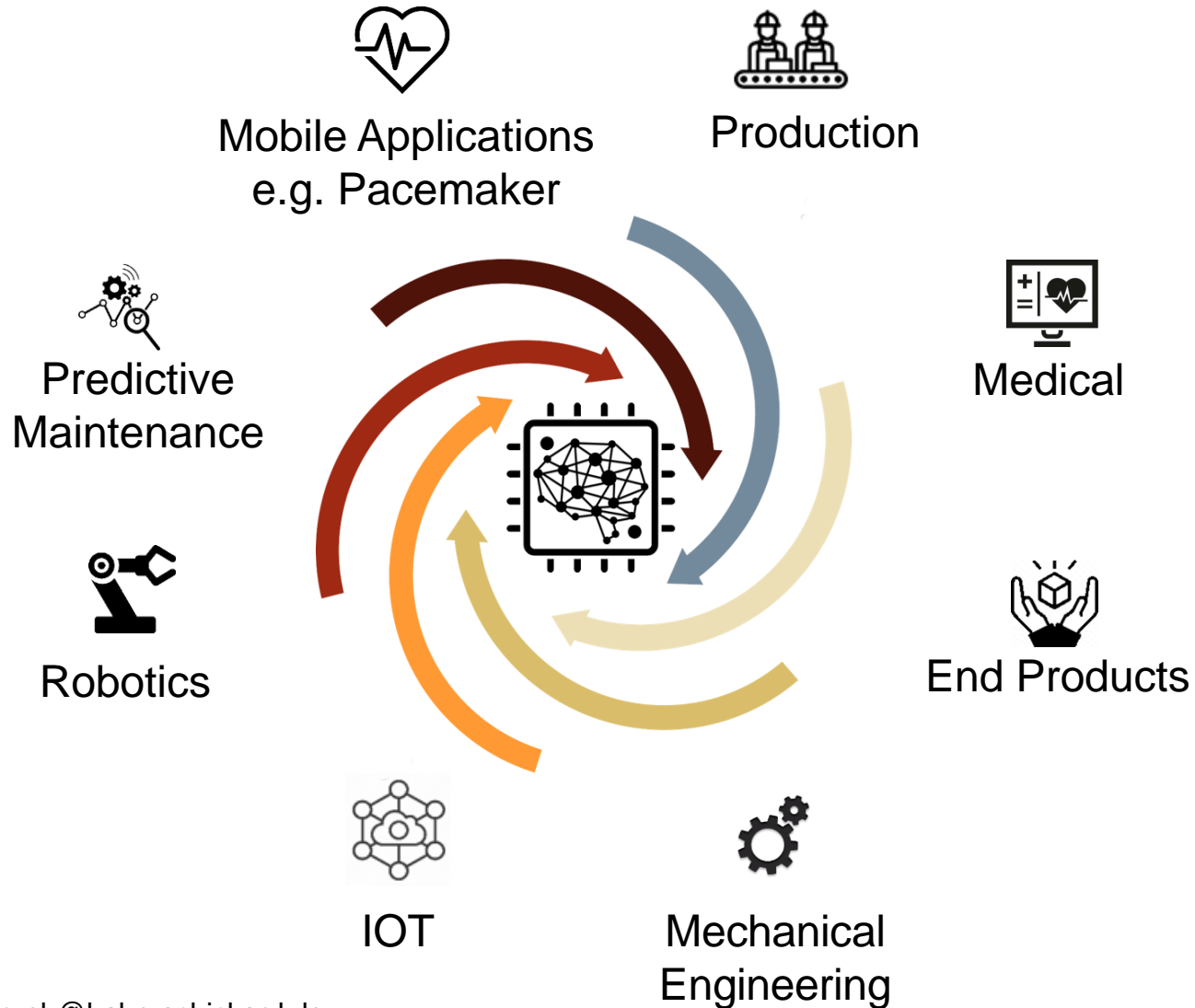
Less Communication



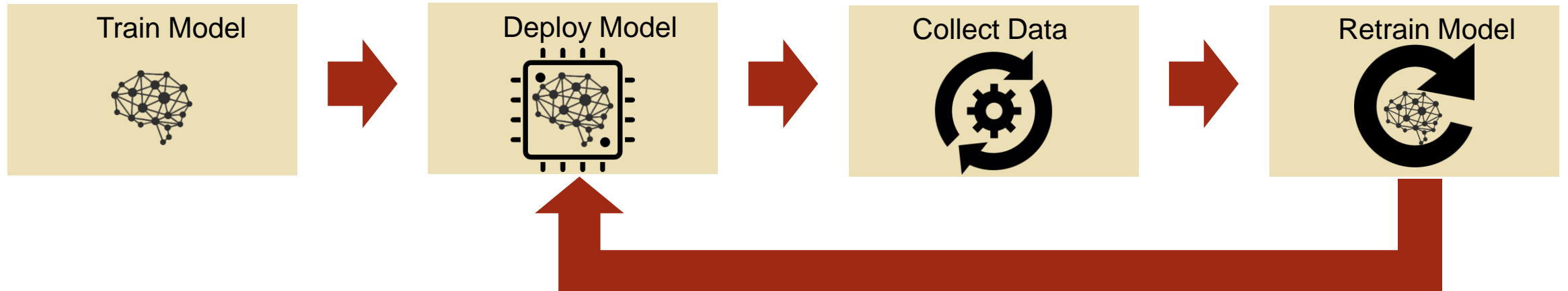
Personalized Models



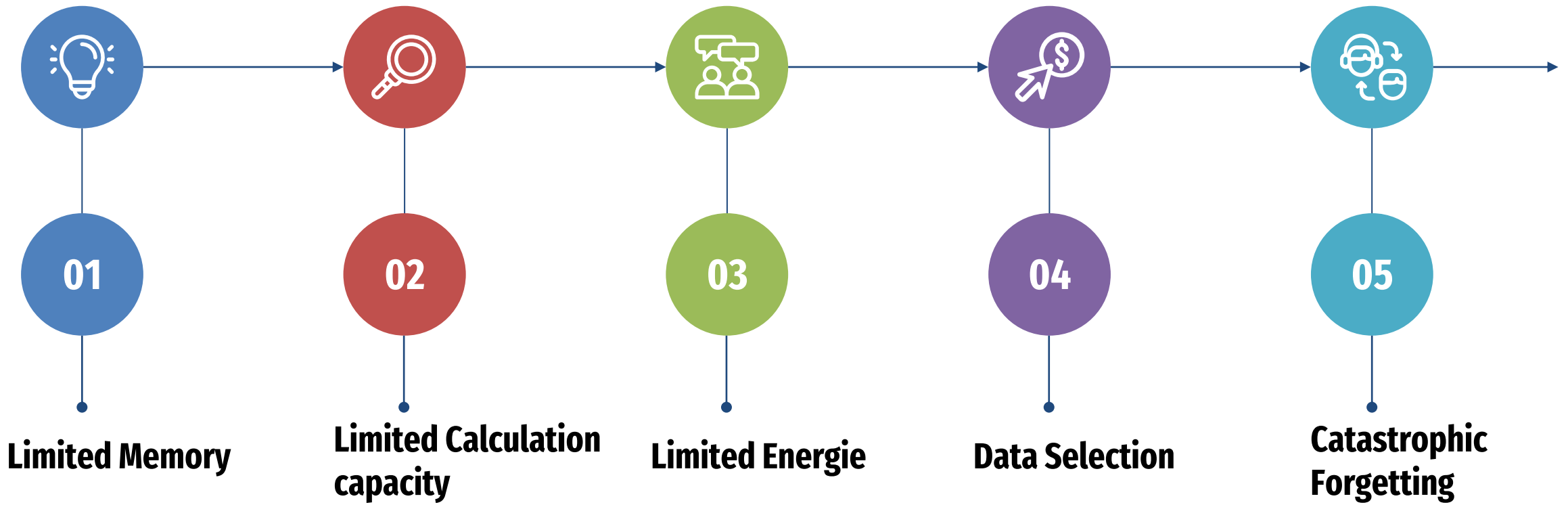
Applications



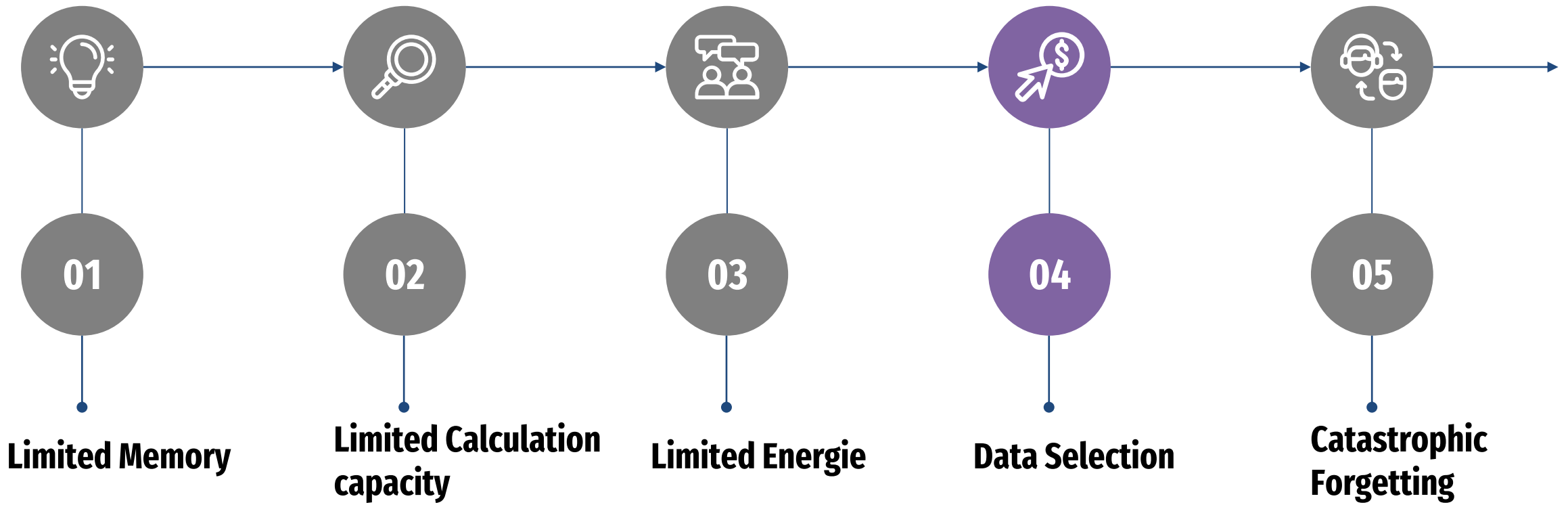
What does the workflow look like?



Problems of OnDevice Training

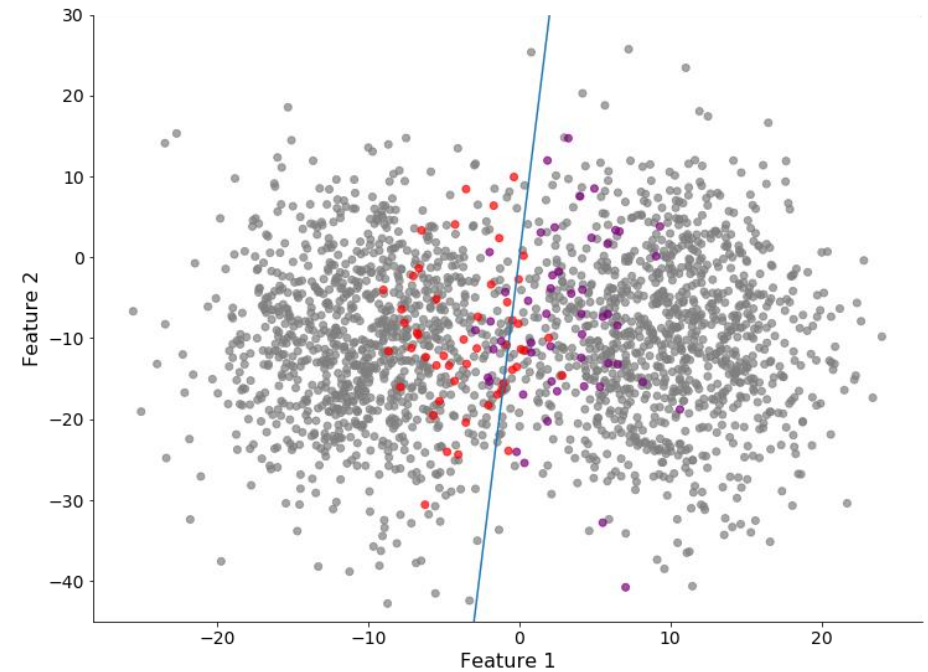
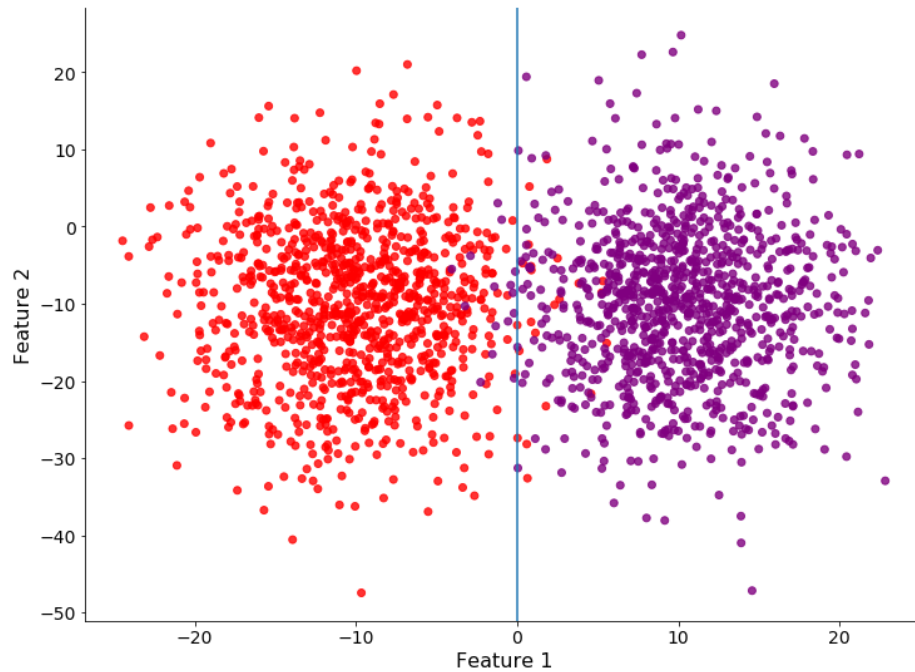


Problems of OnDevice Training

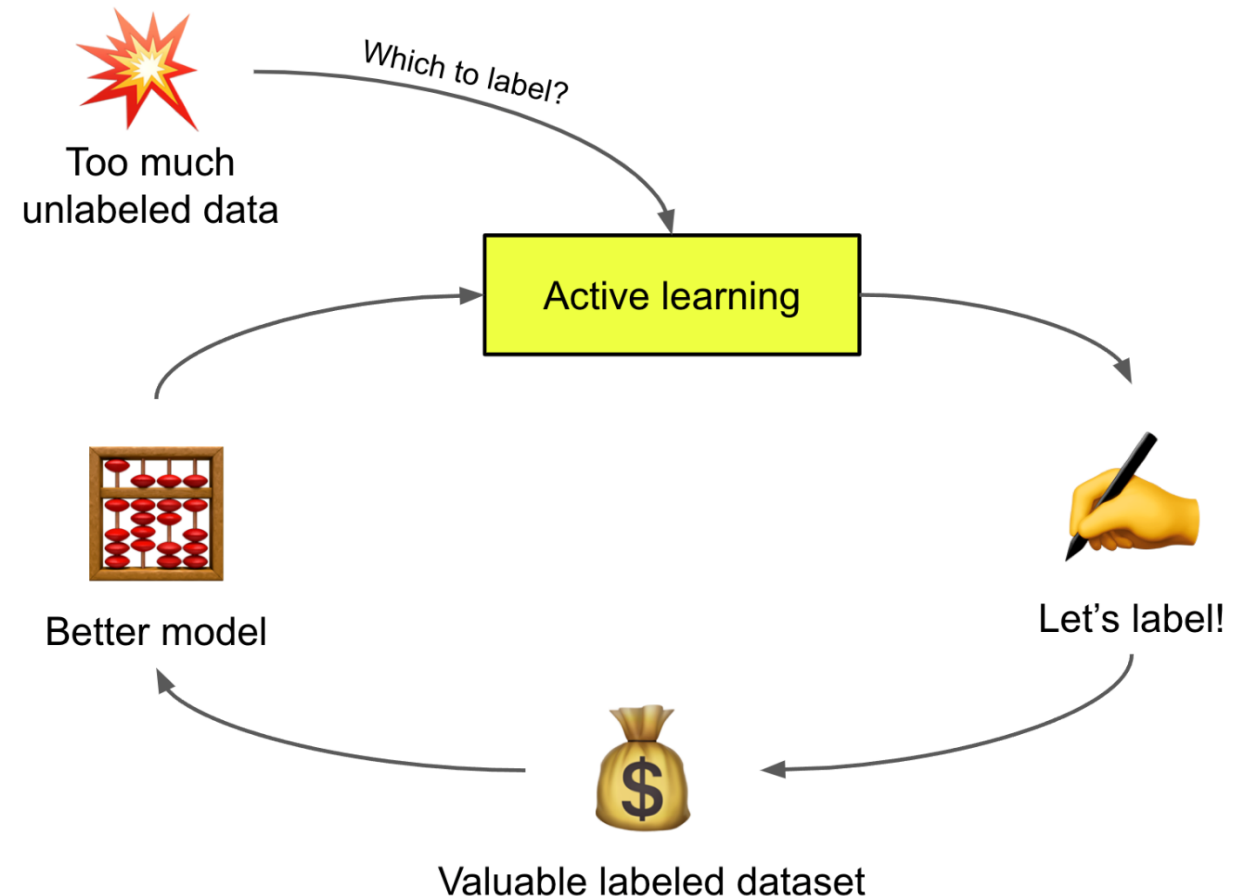


Data Point Selection

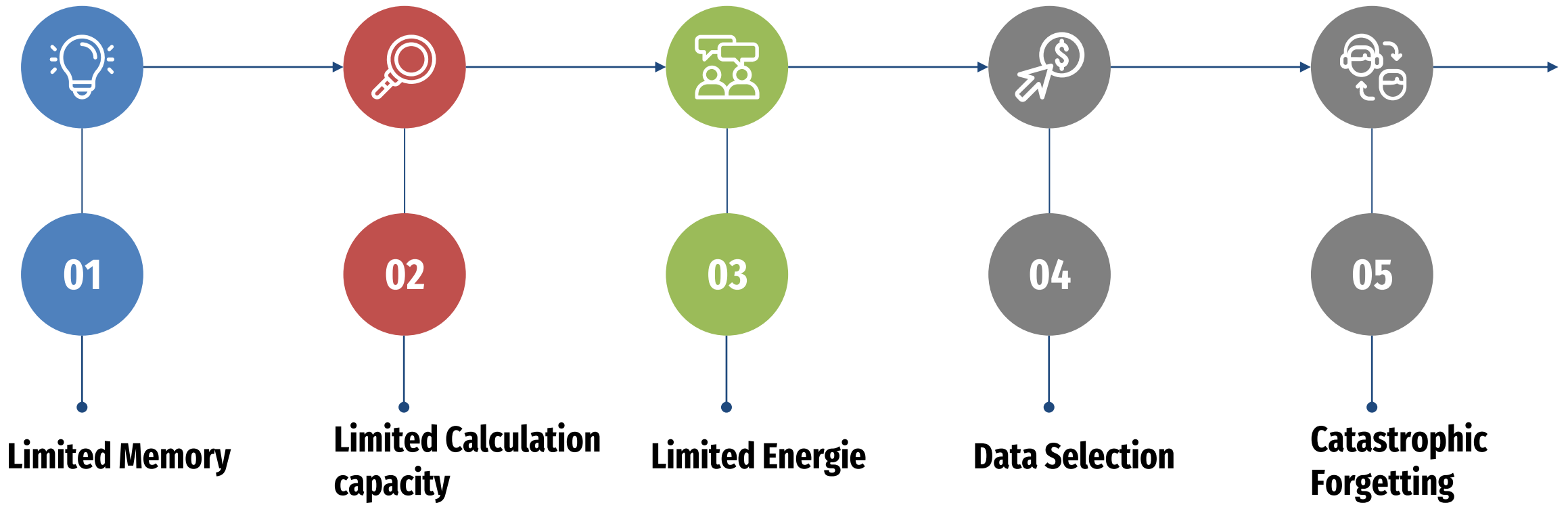
Which data use to train?
How to label Data on Device?
What to do with „old data“? [2]



- 1. Least Confidence:** difference between the most confident prediction and 100% confidence
- 2. Margin of Confidence:** difference between the top two most confident predictions
- 3. Ratio of Confidence:** ratio between the top two most confident predictions
- 4. Entropy:** difference between all predictions, as defined by information theory [3]

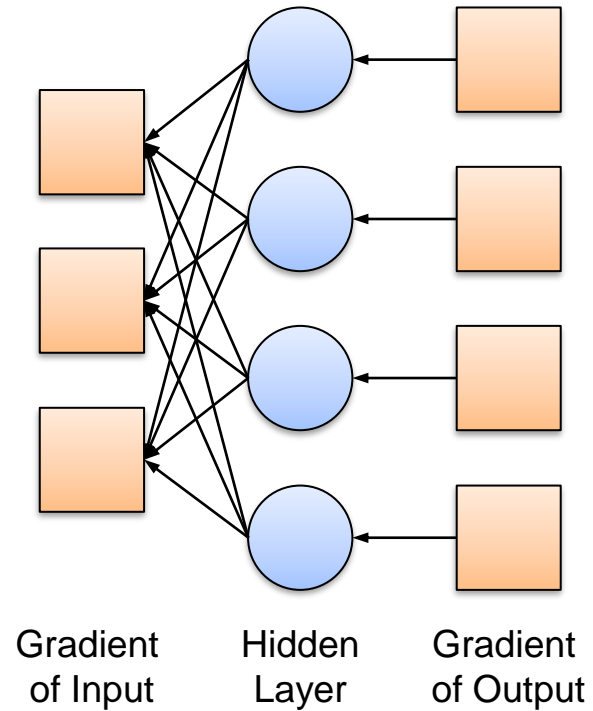


Problems of OnDevice Training

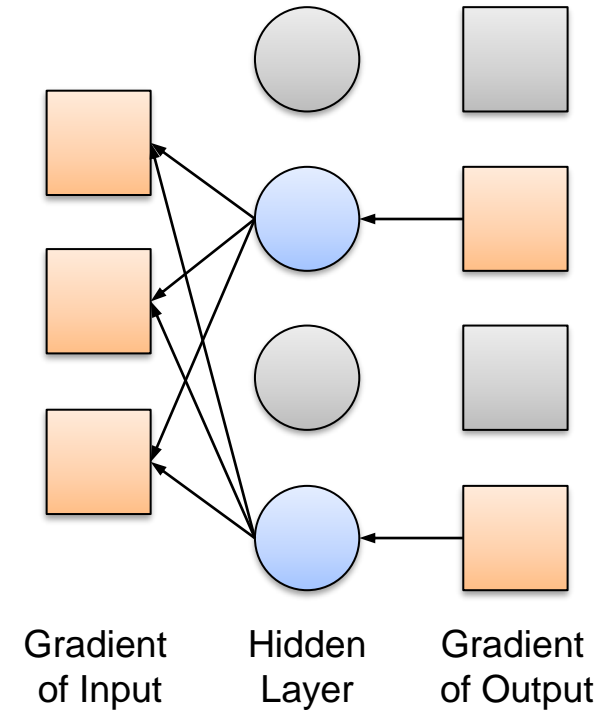


Train more efficient

Backpropagation
100% propagation rate

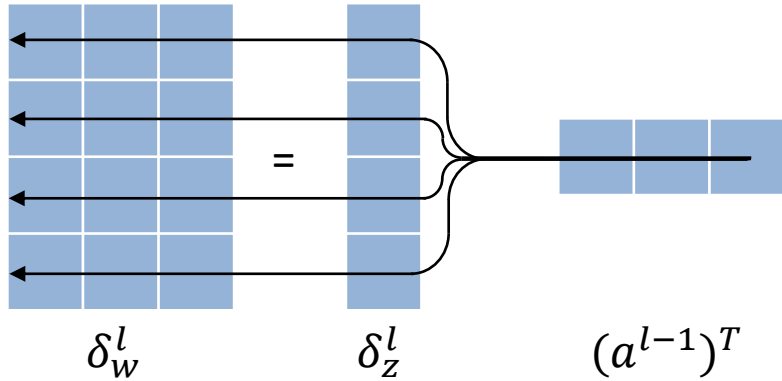


Sparse Backpropagation [4]
50% propagation rate

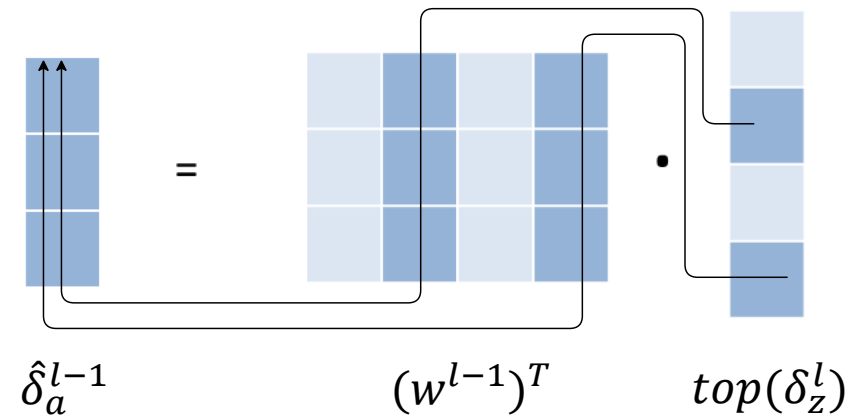
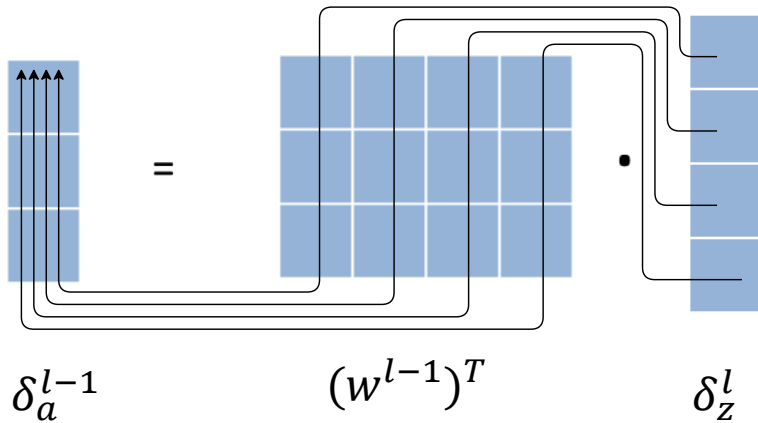
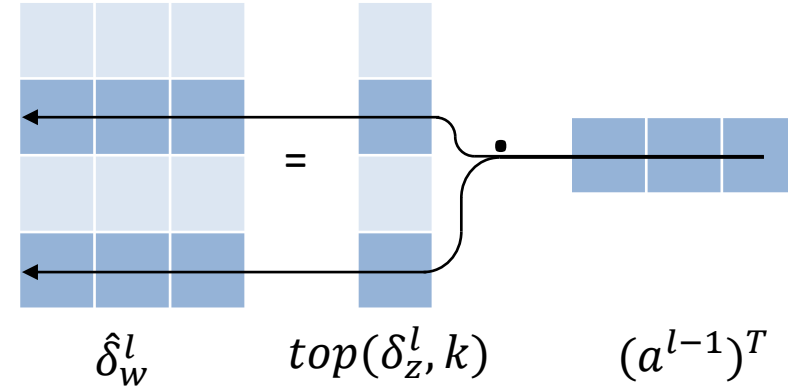


Train more efficient

Backpropagation
100% propagation rate



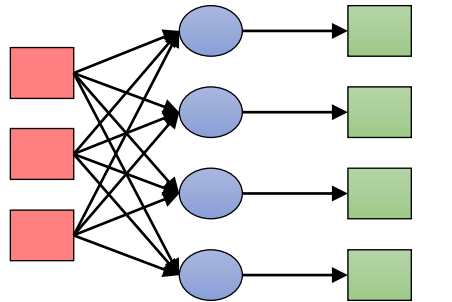
Sparse Backpropagation
50% propagation rate



1. Forward Propagation

$$z^l = W^l \cdot a^{l-1} + b^l$$

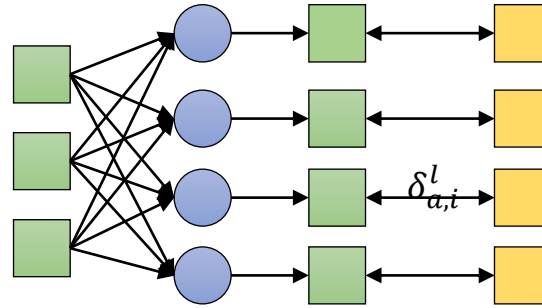
$$a^{l+1} = f(z^l)$$



Input a^{l-1} Hidden Layer Output a^{l+1}

2. Sum magnitudes of local error

$$Y^l = \sum_{i=1}^{N^l} |\delta_{a,i}^l|$$



Input a^{l-1} Hidden Layer Output a^{l+1} Groundtruth

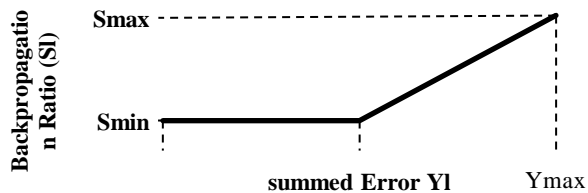
3. Decide to train Datapoint

$$D^L = \left(D_{min} + \alpha^L * \frac{(D_{max} - D_{min})}{\alpha_{max}^L} \right) * \beta^L$$

4. Calculate k

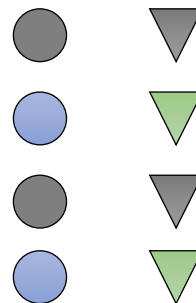
$$S^l := \left(S_{min} + Y^l \cdot \frac{S_{max} - S_{min}}{Y_{max}^l} \right) \cdot \zeta^{L-l}$$

$$k^l = S^l \cdot N^l$$



5. Get Top k

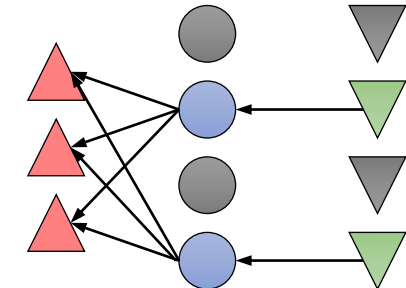
(Top k = 2)



Hidden Layer Gradient of Output

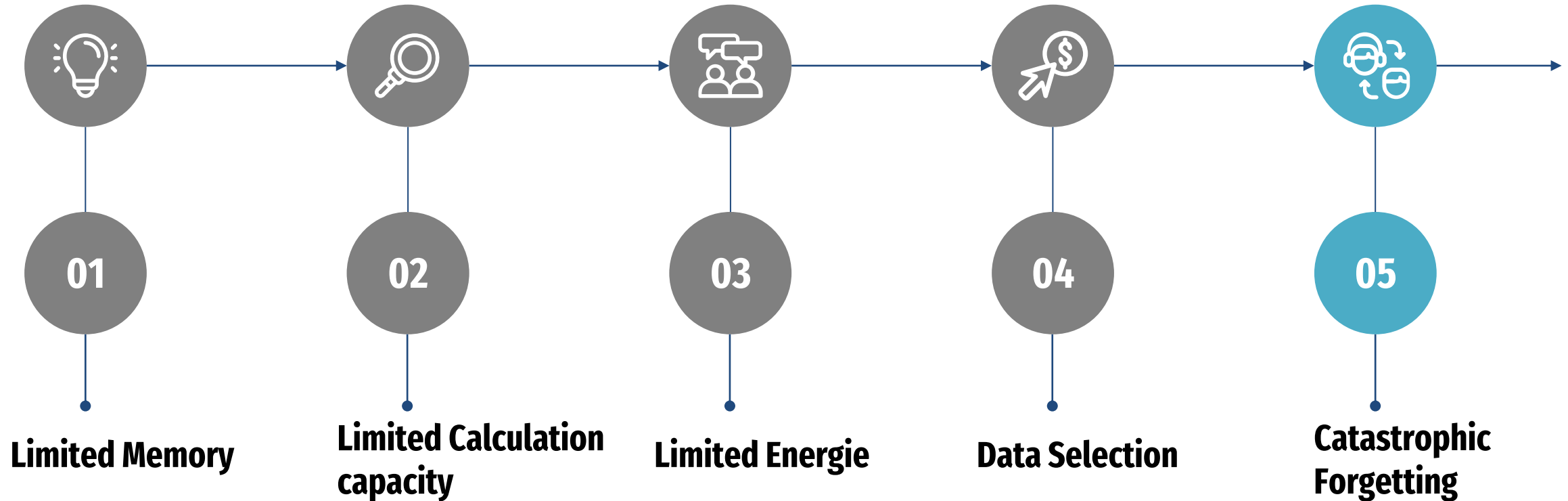
6. Sparse back Propagation

(Top k = 2)

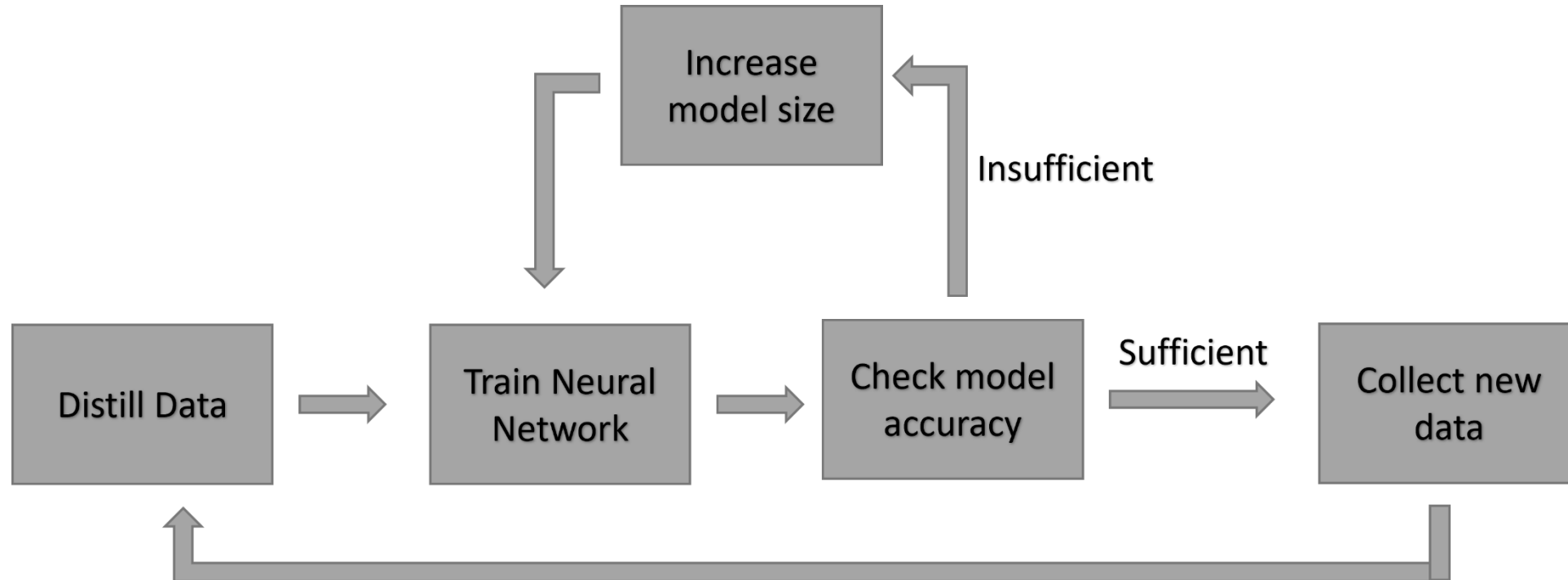


Gradient of Input Hidden Layer Gradient of Output

Problems of OnDevice Training

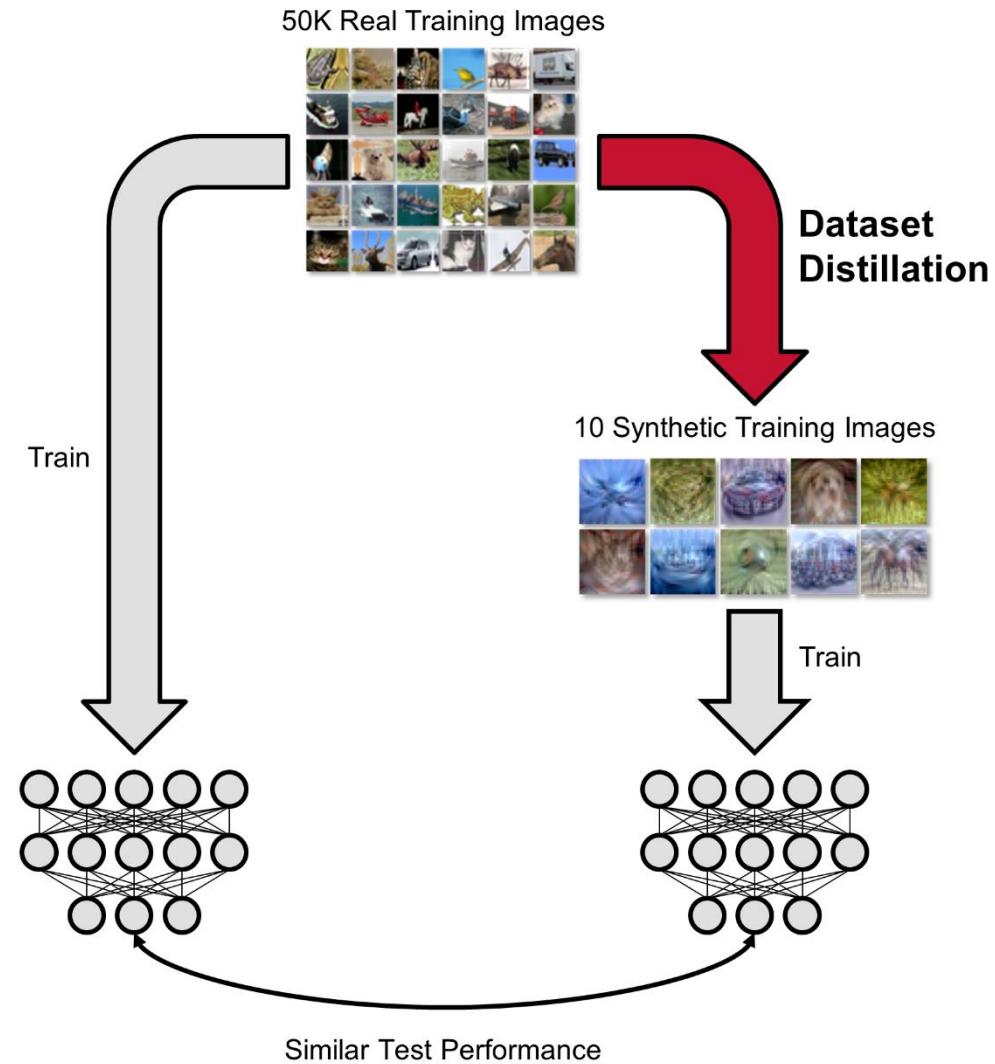


Proposed Workflow



[5]

What is Dataset distillation



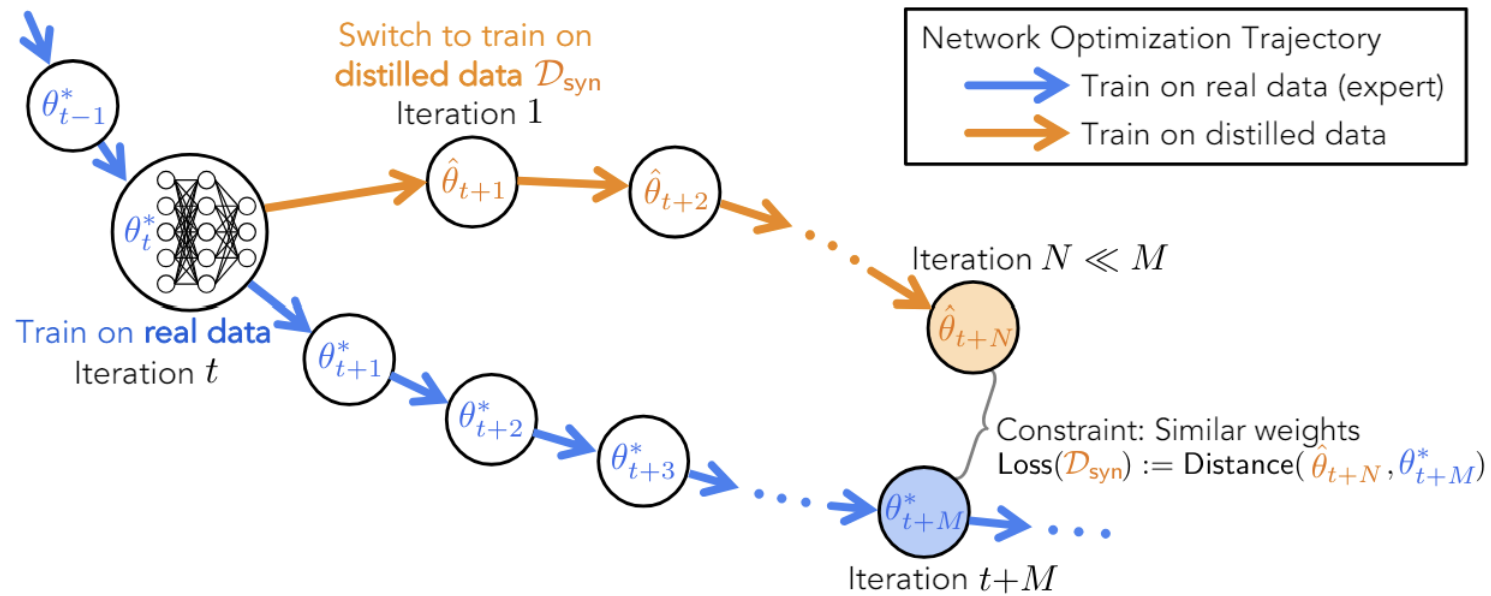
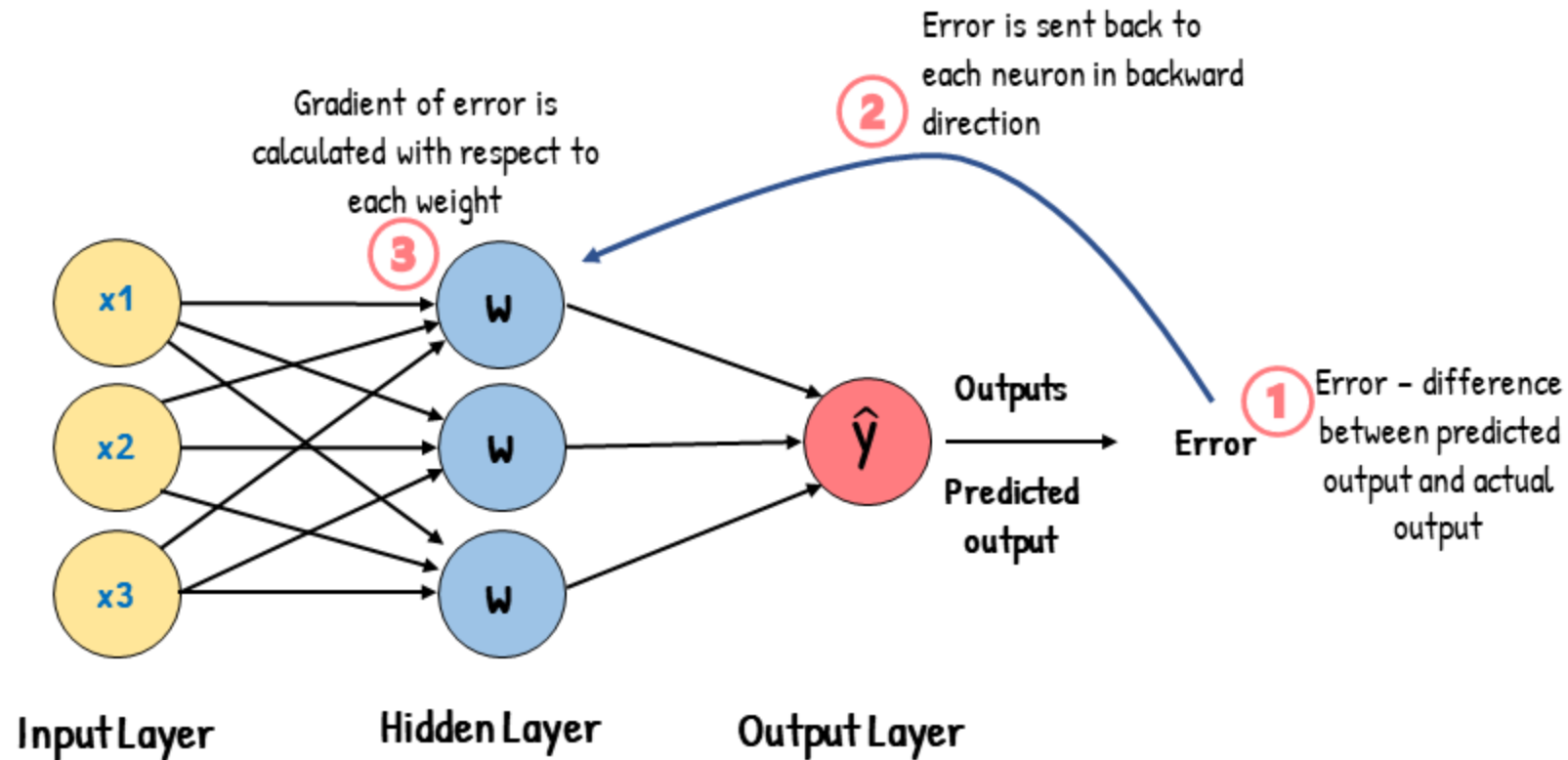


Figure 3. We perform long-range parameter matching between training on distilled synthetic data and training on real data. Starting from the same initial parameters, we train distilled data \mathcal{D}_{syn} such that N training steps on them match the same result (in parameter space) from much more M steps on real data.

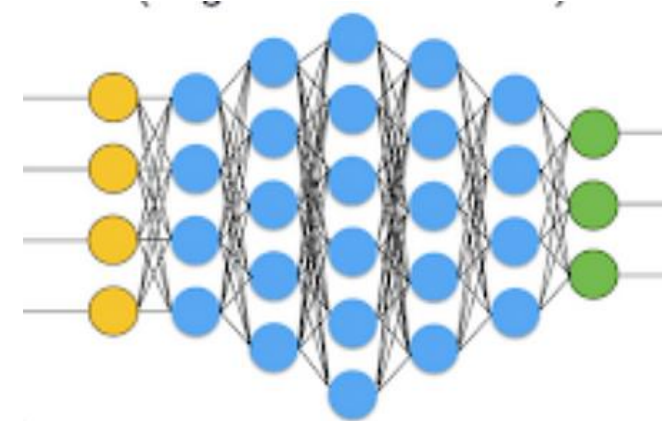
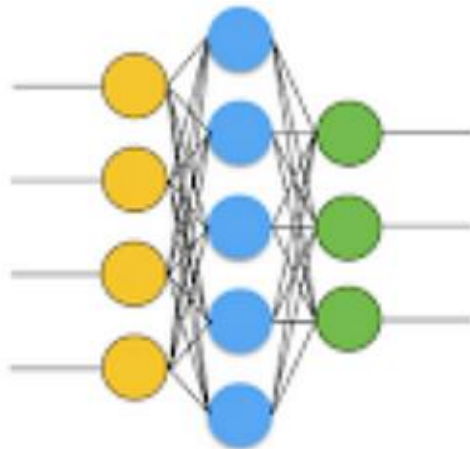
Backpropagation



Check Model Accuracy & Increase Model size

Model accuracy is typically measured using an accuracy score, which is defined as the proportion of correct predictions out of the total number of predictions. If we denote the number of correct predictions as C_{correct} and the total number of predictions as N_{total} , the accuracy score A can be calculated as follows:

$$A = \frac{C_{\text{correct}}}{N_{\text{total}}}$$



Let's denote our original model's architecture as \mathcal{M} and its accuracy score as A . If A doesn't meet the predefined accuracy standard, A_{standard} , we proceed with the enlargement of the model, creating a new model architecture, \mathcal{M}' . This adjustment can be represented as follows:

$$\mathcal{M}' = \begin{cases} \mathcal{M} + \Delta\mathcal{M} & \text{if } A < A_{\text{standard}} \\ \mathcal{M} & \text{otherwise} \end{cases}$$

where $\Delta\mathcal{M}$ represents the increase in the model's complexity. This increase can be in the form of additional layers, more

- [1] Marcus Rüb, Prof. Dr. Axel Sikora, A Practical View on Training Neural Networks in the Edge, IFAC-PapersOnLine, Volume 55, Issue 4, 2022, Pages 272-279, ISSN 2405-8963
- [2] A Comparative Survey of Deep Active Learning.
Xueying Zhan, Qingzhong Wang, Kuan-hao Huang, Haoyi Xiong, Dejing Dou, and Antoni B. Chan,
arXiv:2203.13450
- [3] <https://towardsdatascience.com/uncertainty-sampling-cheatsheet-ec57bc067c0b>
- [4] Rüb, Marcus & Maier, Daniel & Mueller-Gritschneider, Daniel & Sikora, Axel. (2023).
TinyProp -- Adaptive Sparse Backpropagation for Efficient TinyML On-device Learning.
GITHUB: <https://github.com/r1marcus/TinyProp>
- [5] Rüb, Mueller-Gritschneider, Sikora. (2024).
A Continual and Incremental Learning Approach for TinyML On-device Training Using Dataset Distillation and Model Size Adaption.

Thank you





**Hahn
Schickard**

Marcus Rüb
Wissenschaftlicher Mitarbeiter

**Hahn-Schickard-Gesellschaft
für angewandte Forschung e.V.**

Wilhelm-Schickard-Str. 10,
78052 Villingen-Schwenningen
Telefon: +49 7721 943-180
Marcus.Rueb@Hahn-Schickard.de
www.Hahn-Schickard.de