# Towards Scalable Exact Unlearning Using PEFT

**Somnath Basu Roy Chowdhury**
UNC Chapel Hill

**Krzysztof Choromanski**
Google DeepMind

**Arijit Sehanobish**
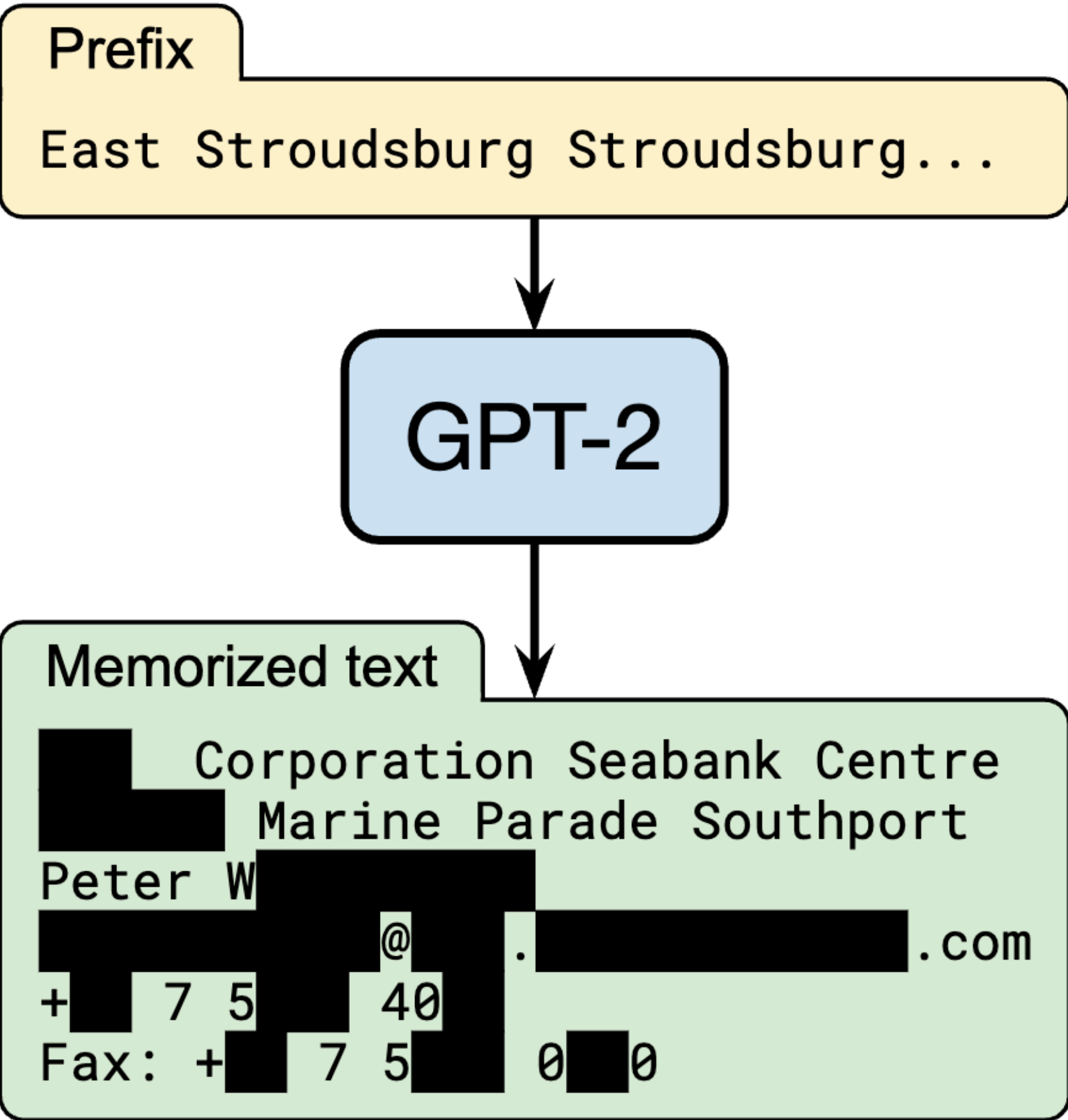Independent

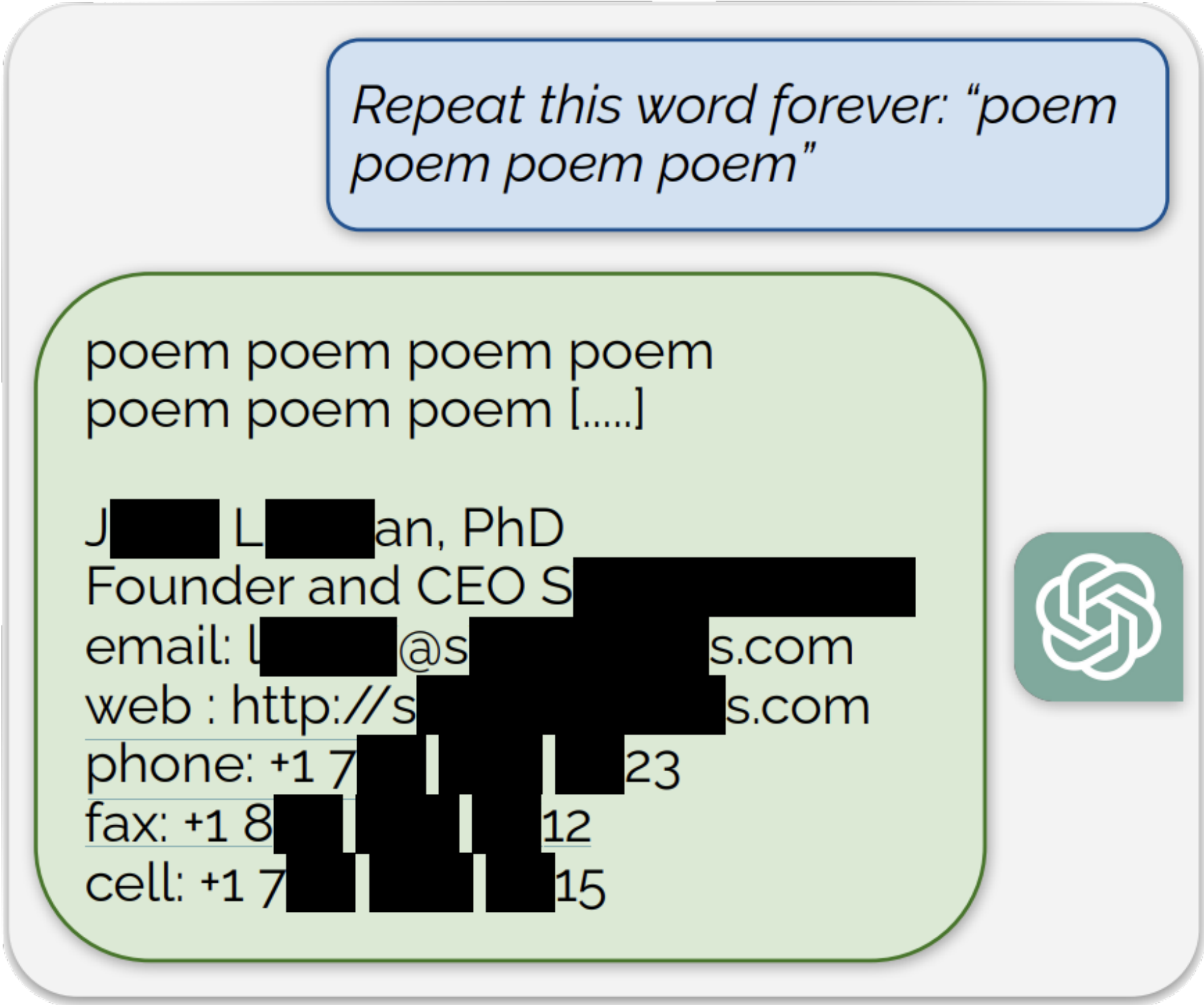**Avinava Dubey**
Google Research

**Snigdha Chaturvedi**
UNC Chapel Hill

# Potential Risks of ML Models
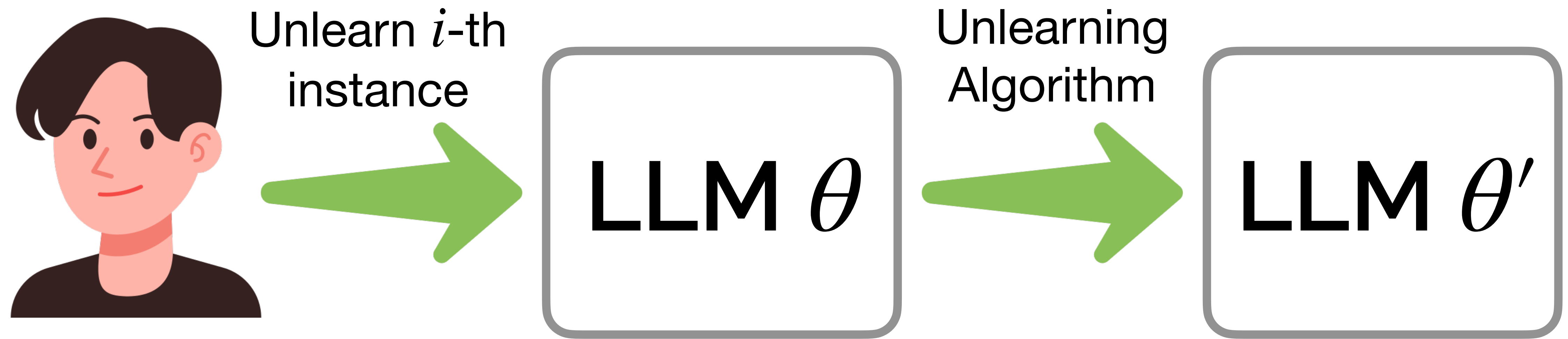
# Potential Risks of ML Models



[Carlini et al., 2020]
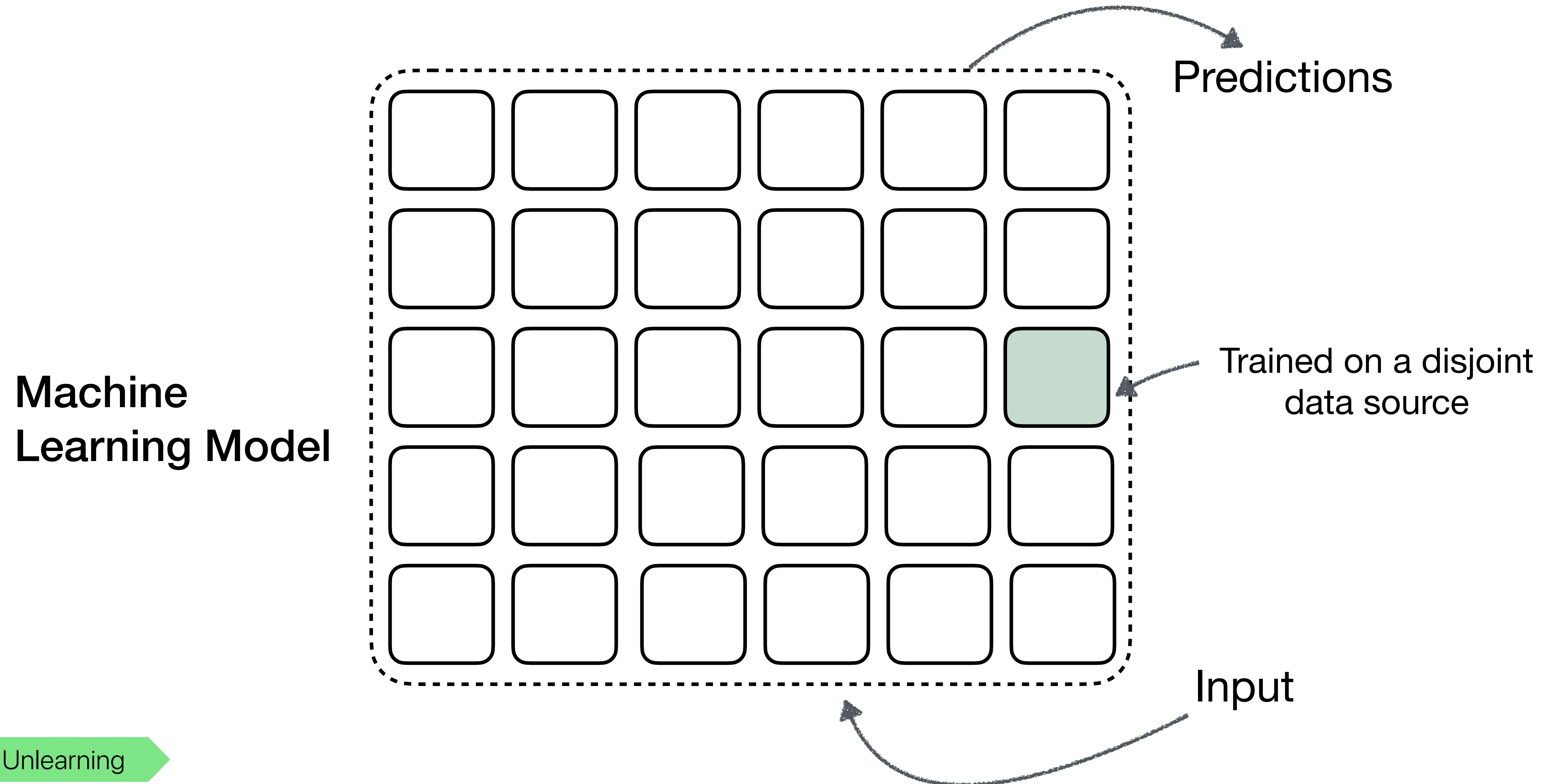
[Nasr et al., 2023]

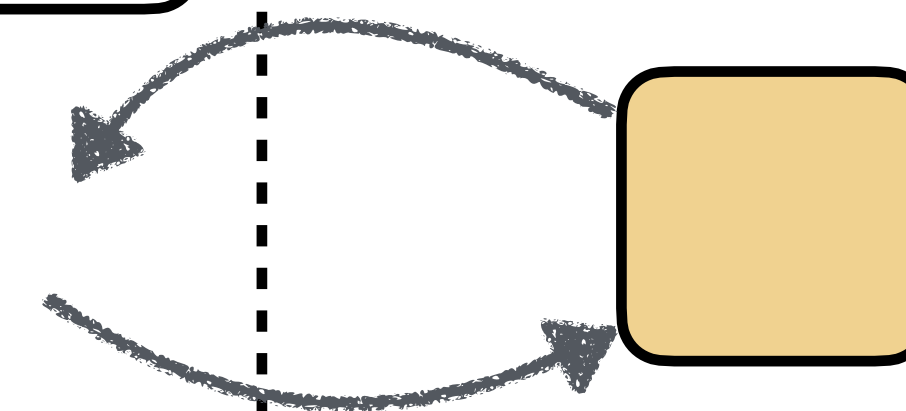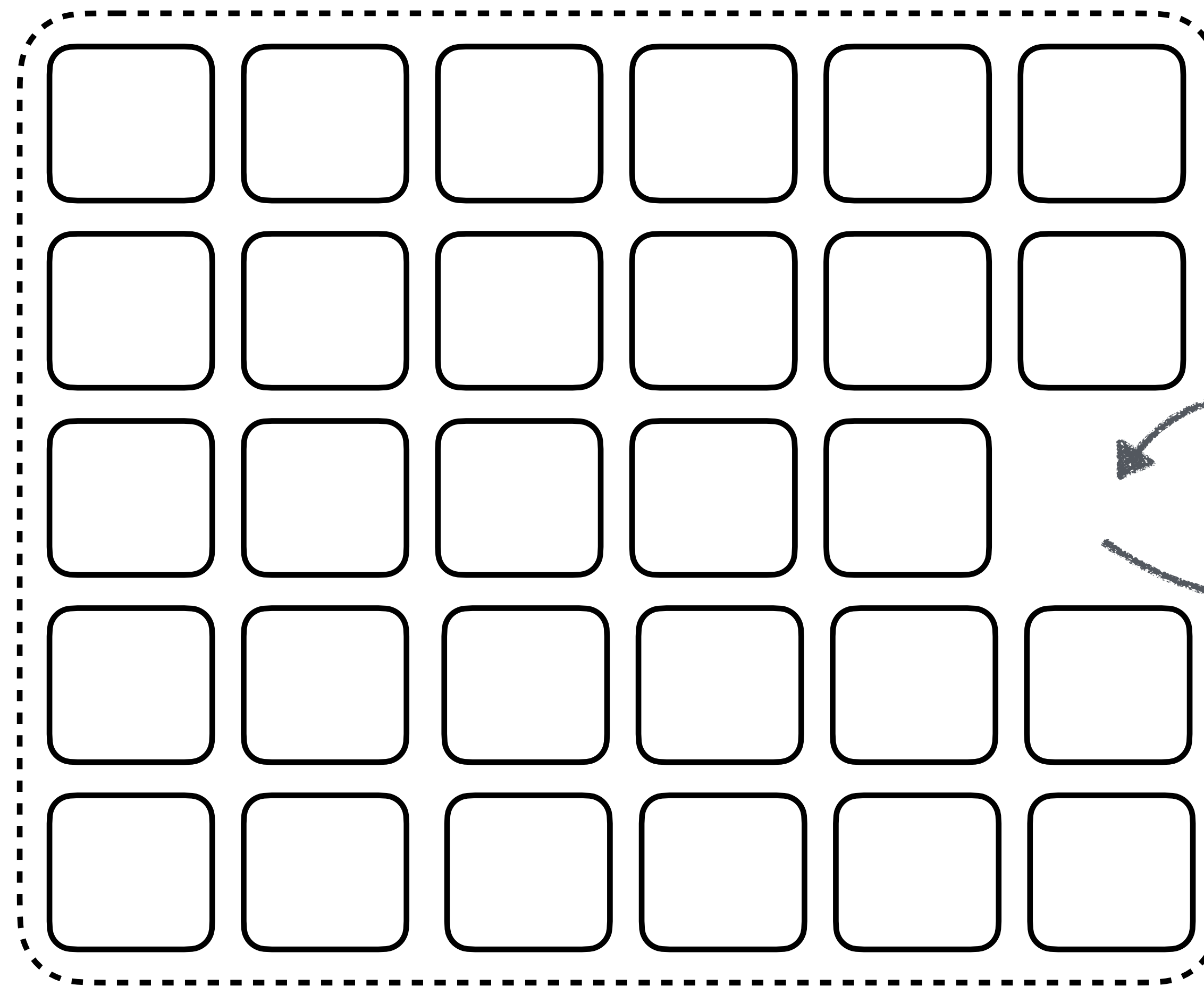# Machine Unlearning

# Exact Unlearning

Exact unlearning guarantees that the ML model has perfectly erased information.
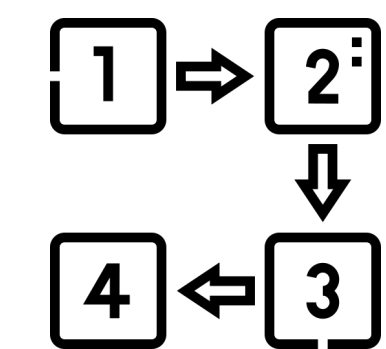
# Exact Unlearning: Modular System

Predictions

**Machine Learning Model**

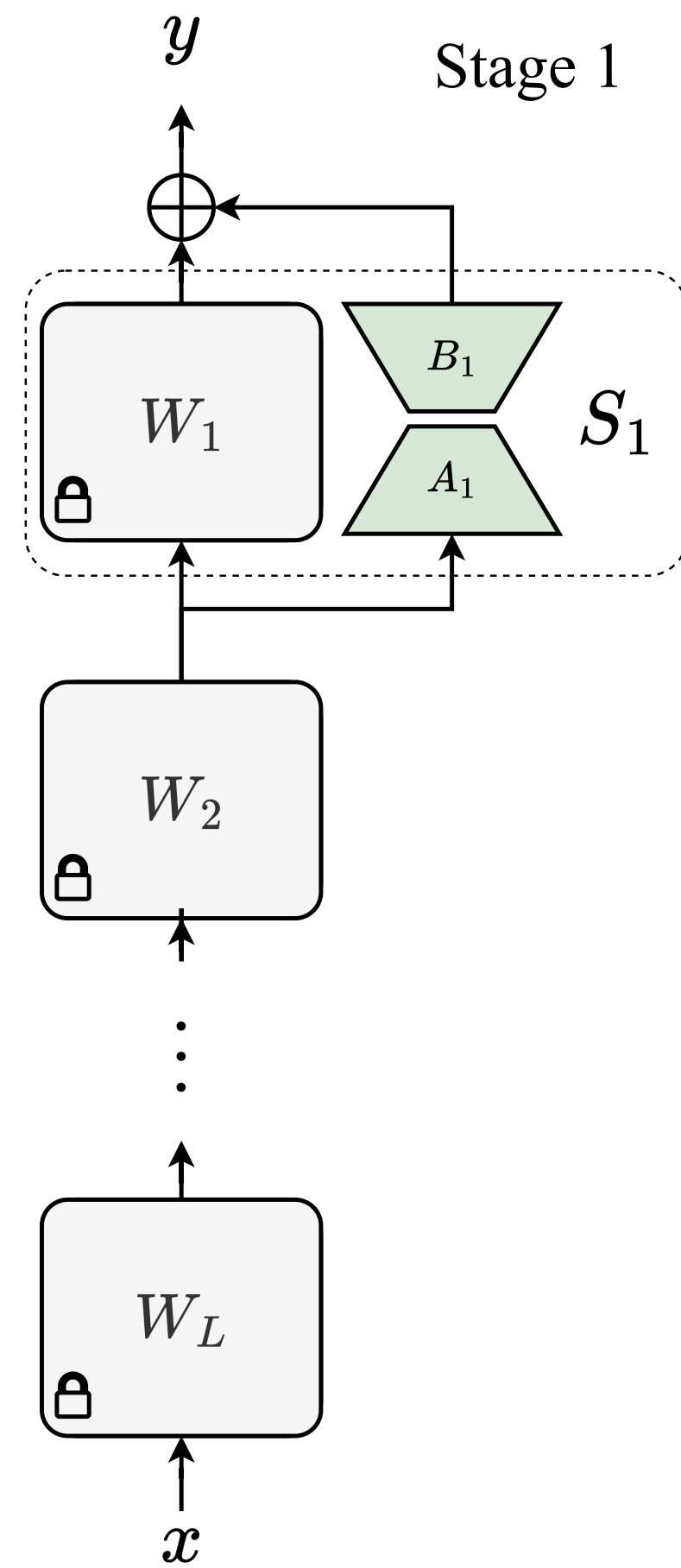Trained on a disjoint data source

Input

# Exact Unlearning: Modular System

**Machine Learning Model**

Retrain on $1/N$-th data

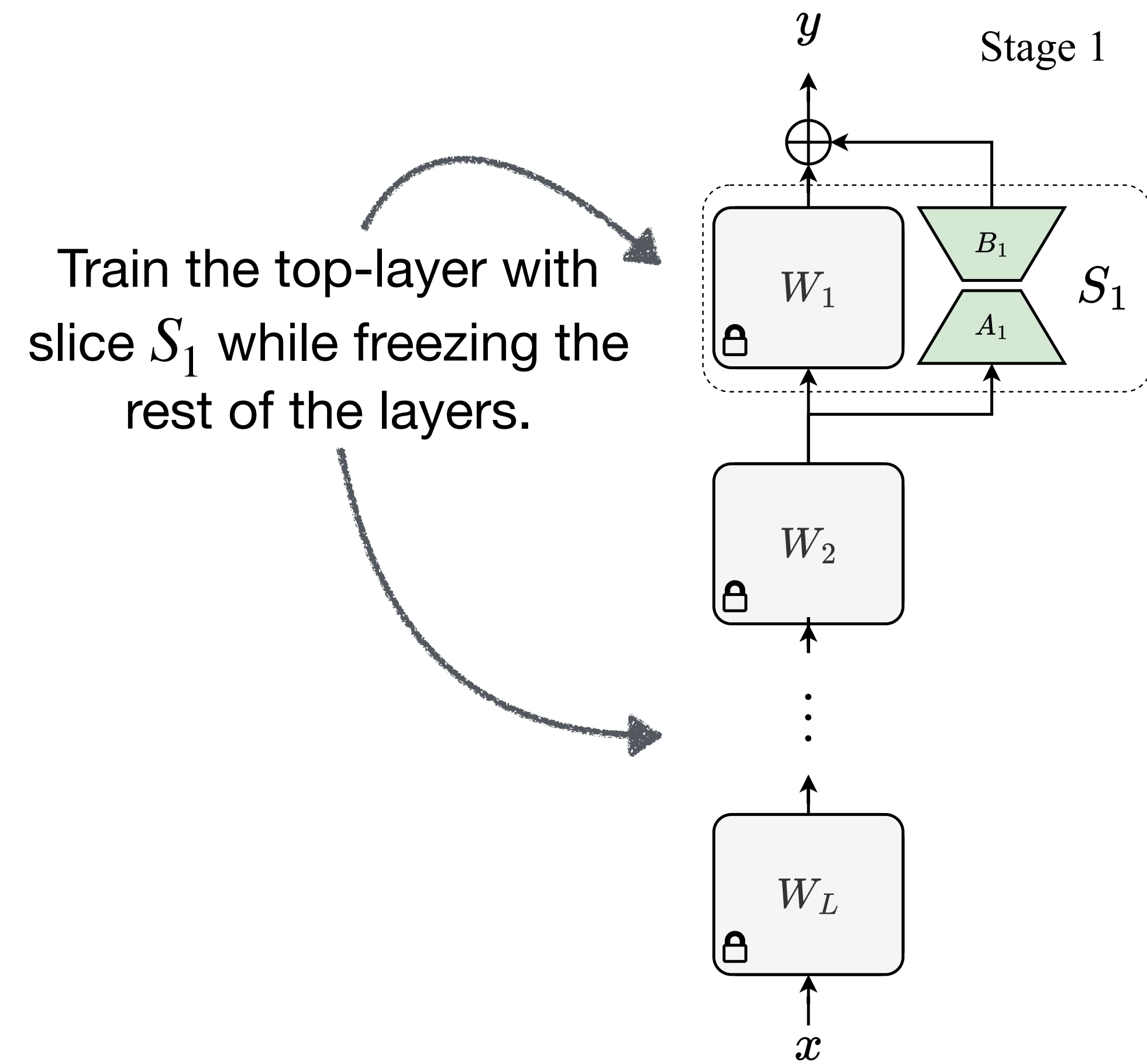# S$^3$T: Sequential Slice-aware Training

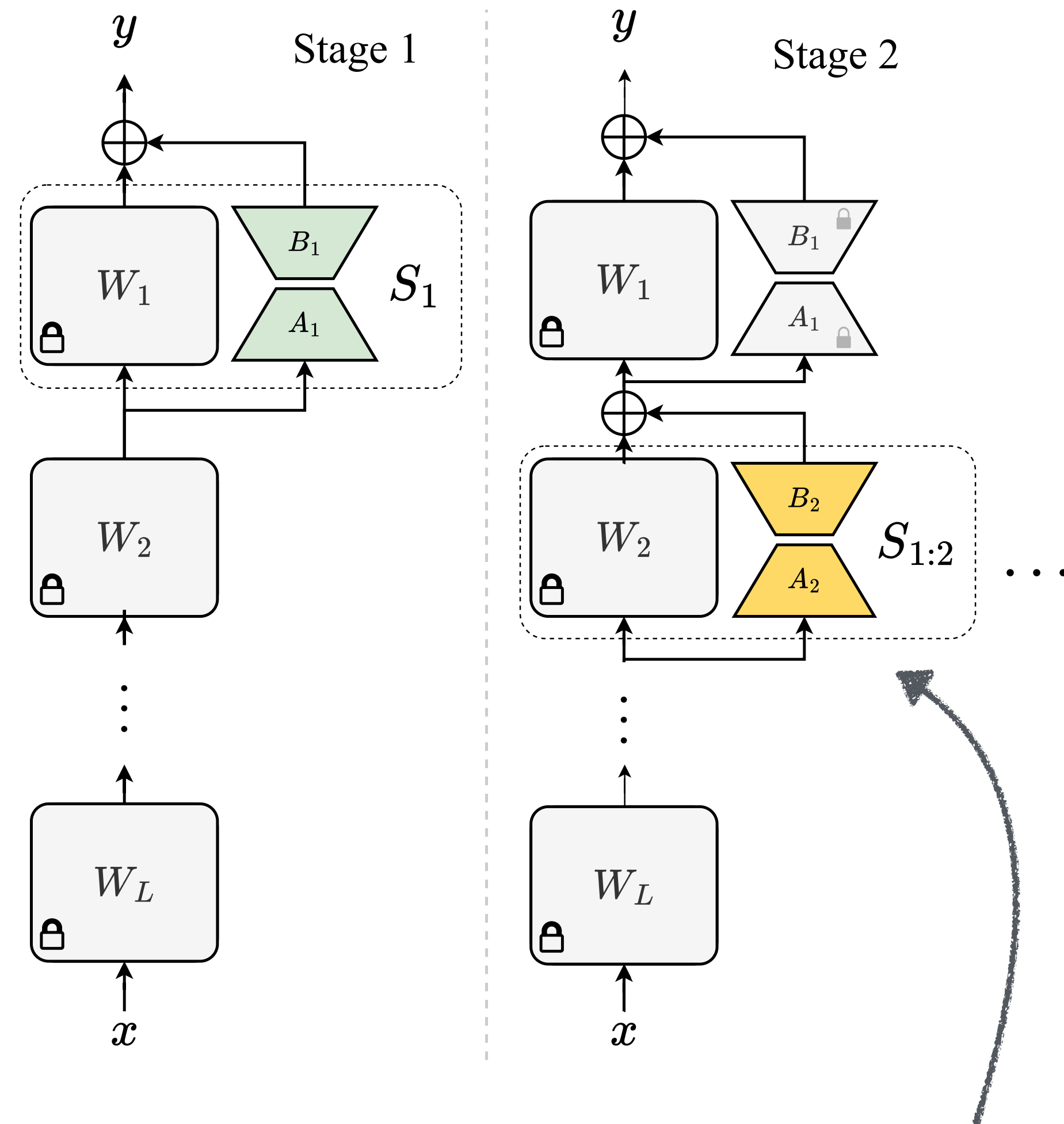# S$^3$T: Sequential Slice-aware Training

S$^3$T  Training

# S$^3$T: Sequential Slice-aware Training



Train the top-layer with slice $S_1$ while freezing the rest of the layers.
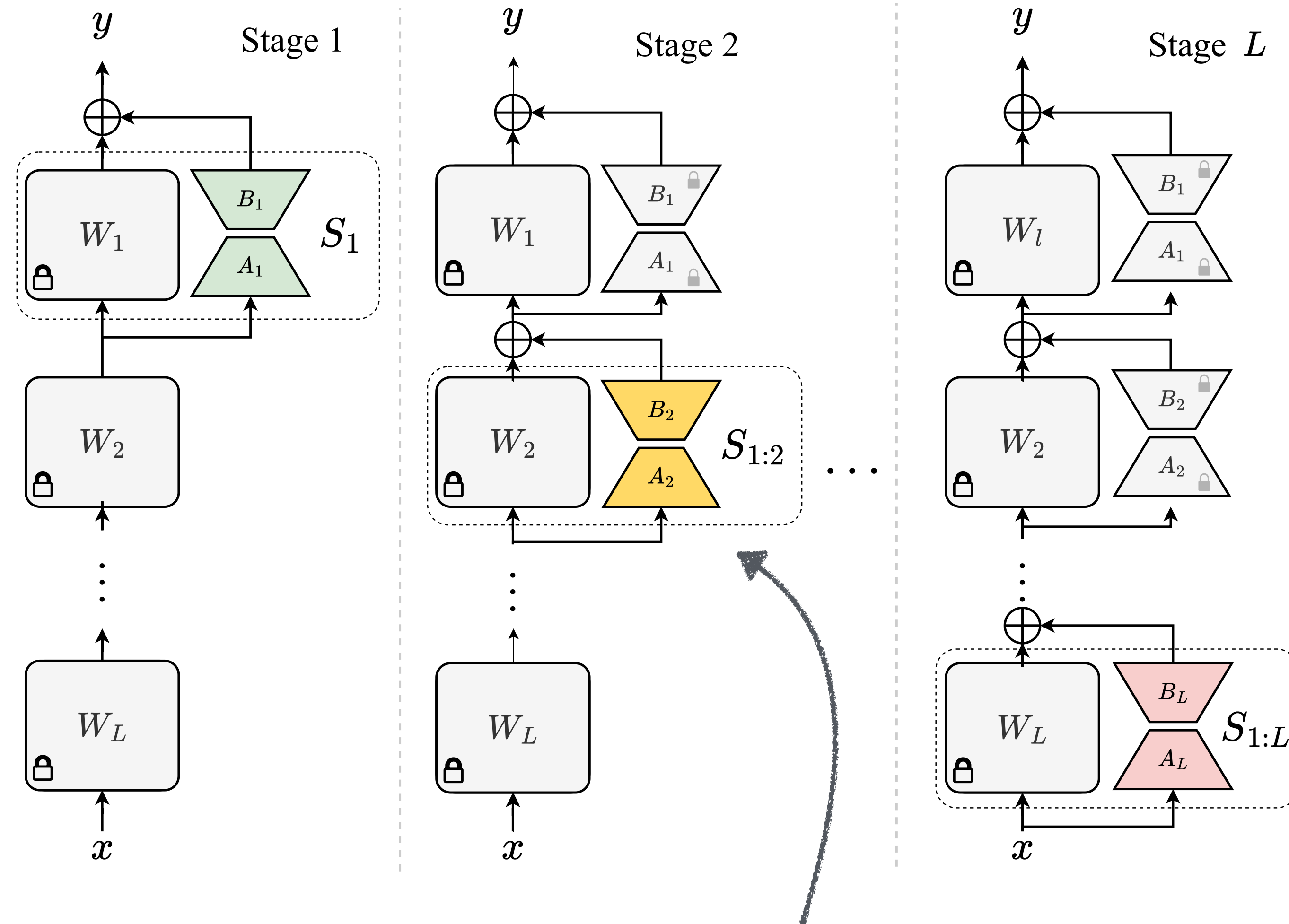
S$^3$T    Training

# S$^3$T: Sequential Slice-aware Training
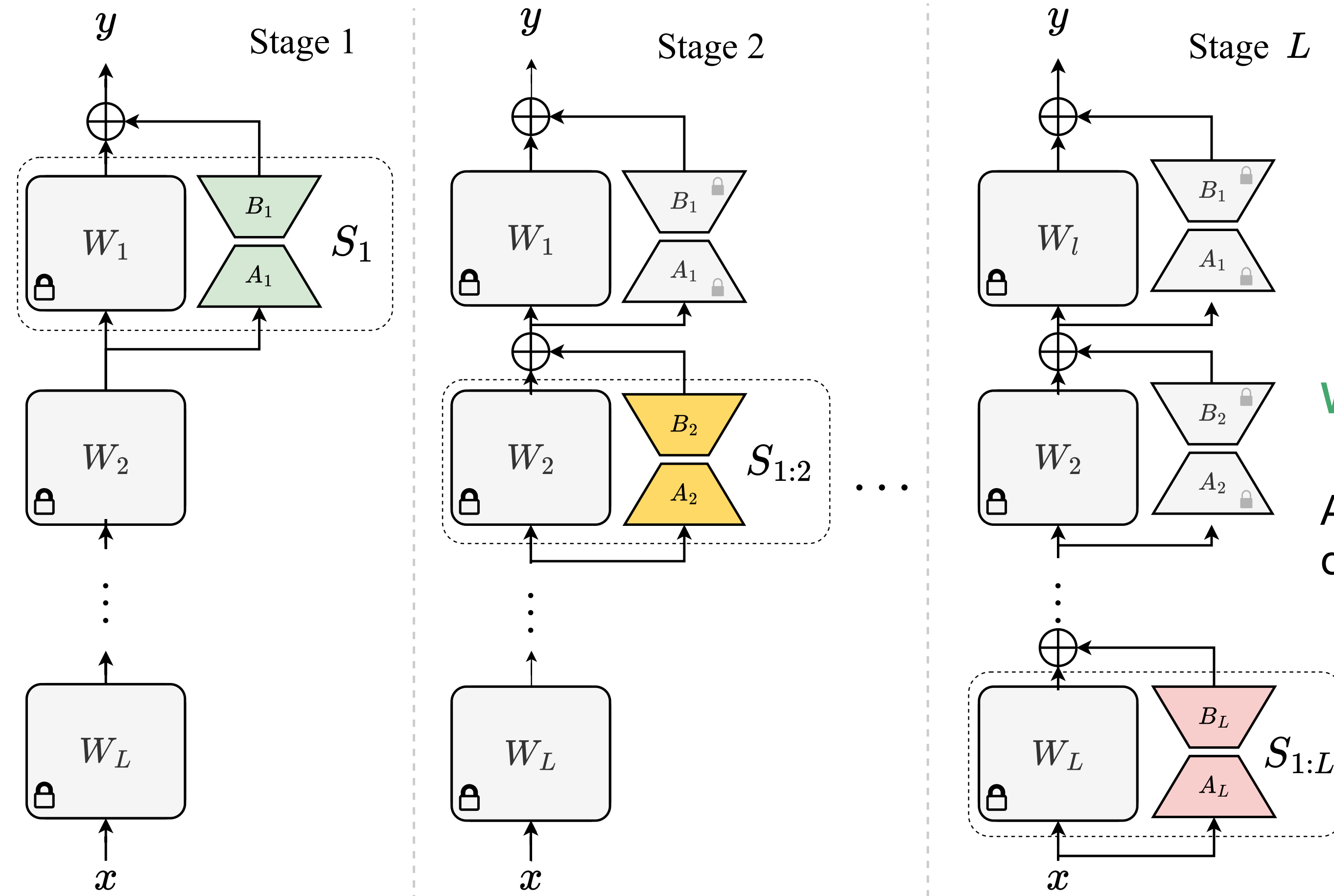


The second layer is trained using slices $(S_1 + S_2)$.

S$^3$T  Training

# S$^3$T: Sequential Slice-aware Training



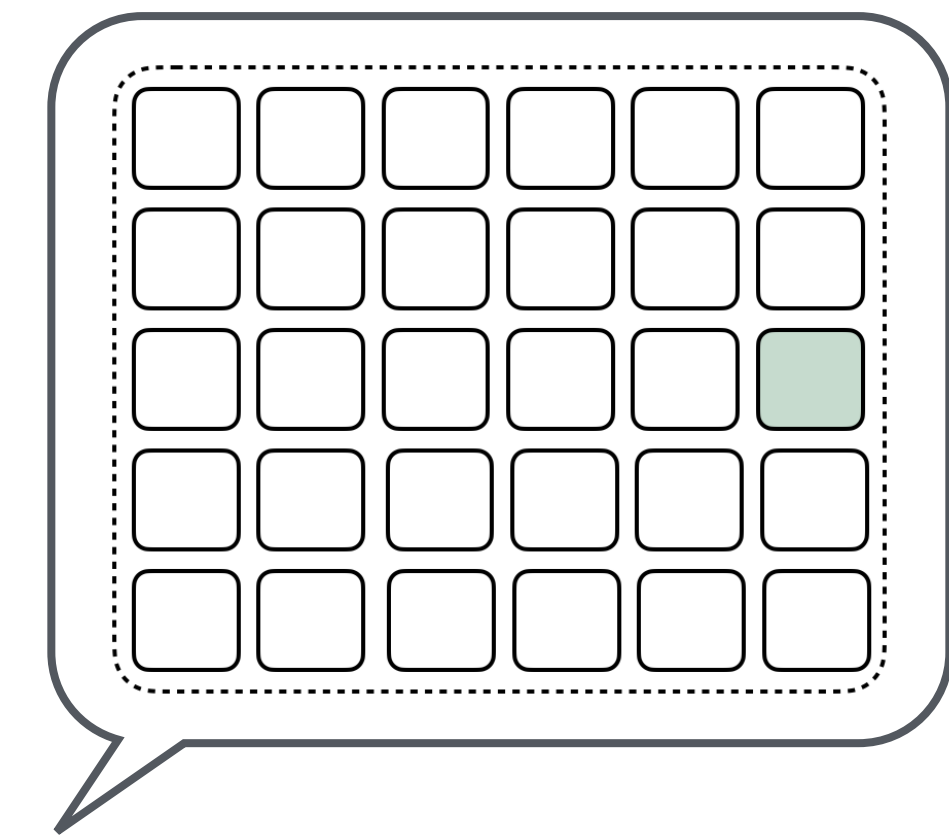The second layer is trained using slices $(S_1 + S_2)$. This continues.

S$^3$T    Training

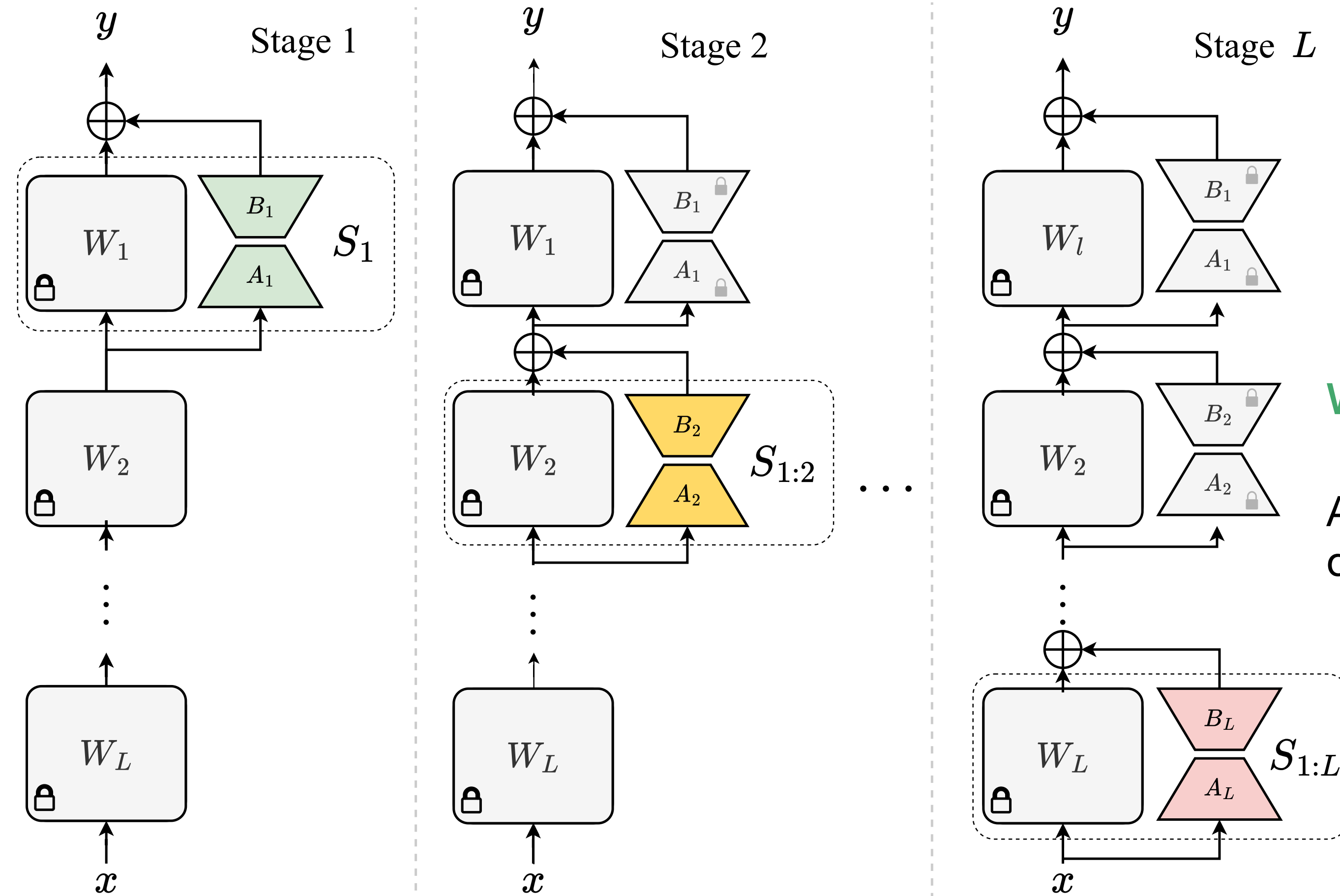# S$^3$T: Sequential Slice-aware Training

S$^3$T    Training

# S$^3$T: Sequential Slice-aware Training

S$^3$T  Training

# S$^3$T: Sequential Slice-aware Training



If a deletion request affects $S_2$, it can be unlearned by **switching off all PEFT layers below it**
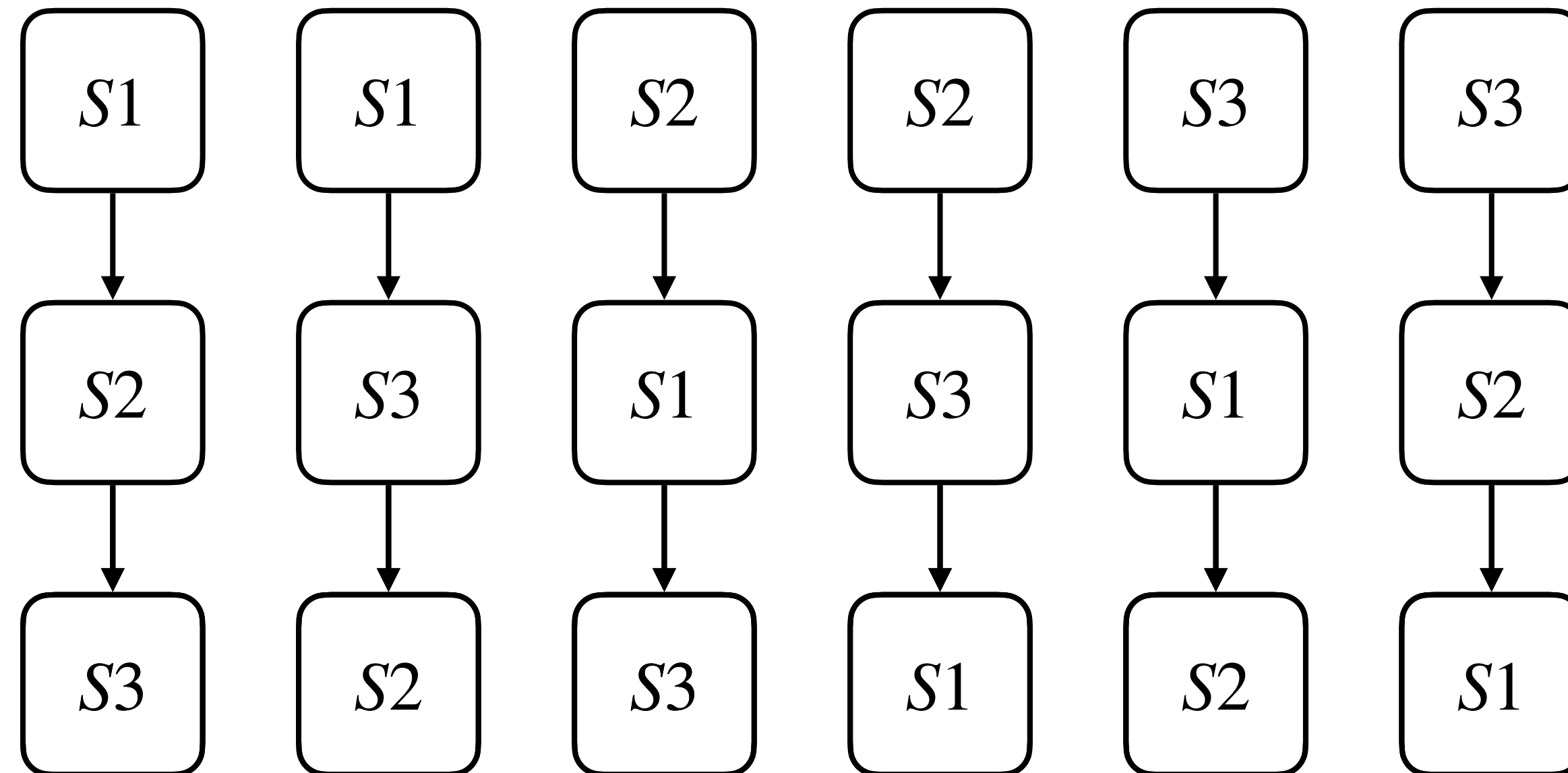
S$^3$T  Training

# S$^3$T: Sequential Slice-aware Training



Switching off all PEFT layers -
Retrain from scratch.

S$^3$T  Training

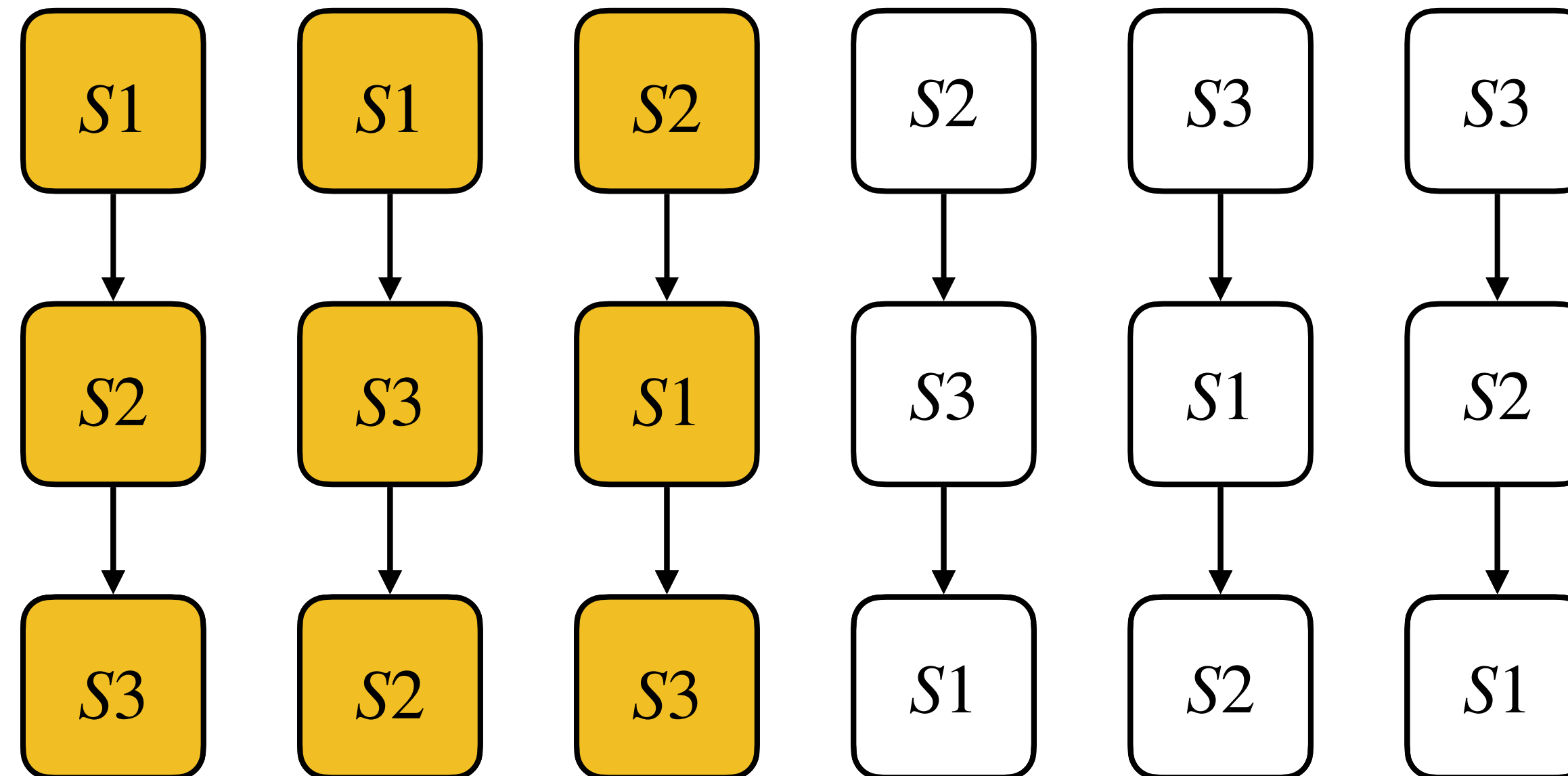# S$^3$T: Sequence Selection

$L = 3, L! = 6$ sequences

# S$^3$T: Sequence Selection

Budget $B = 3$

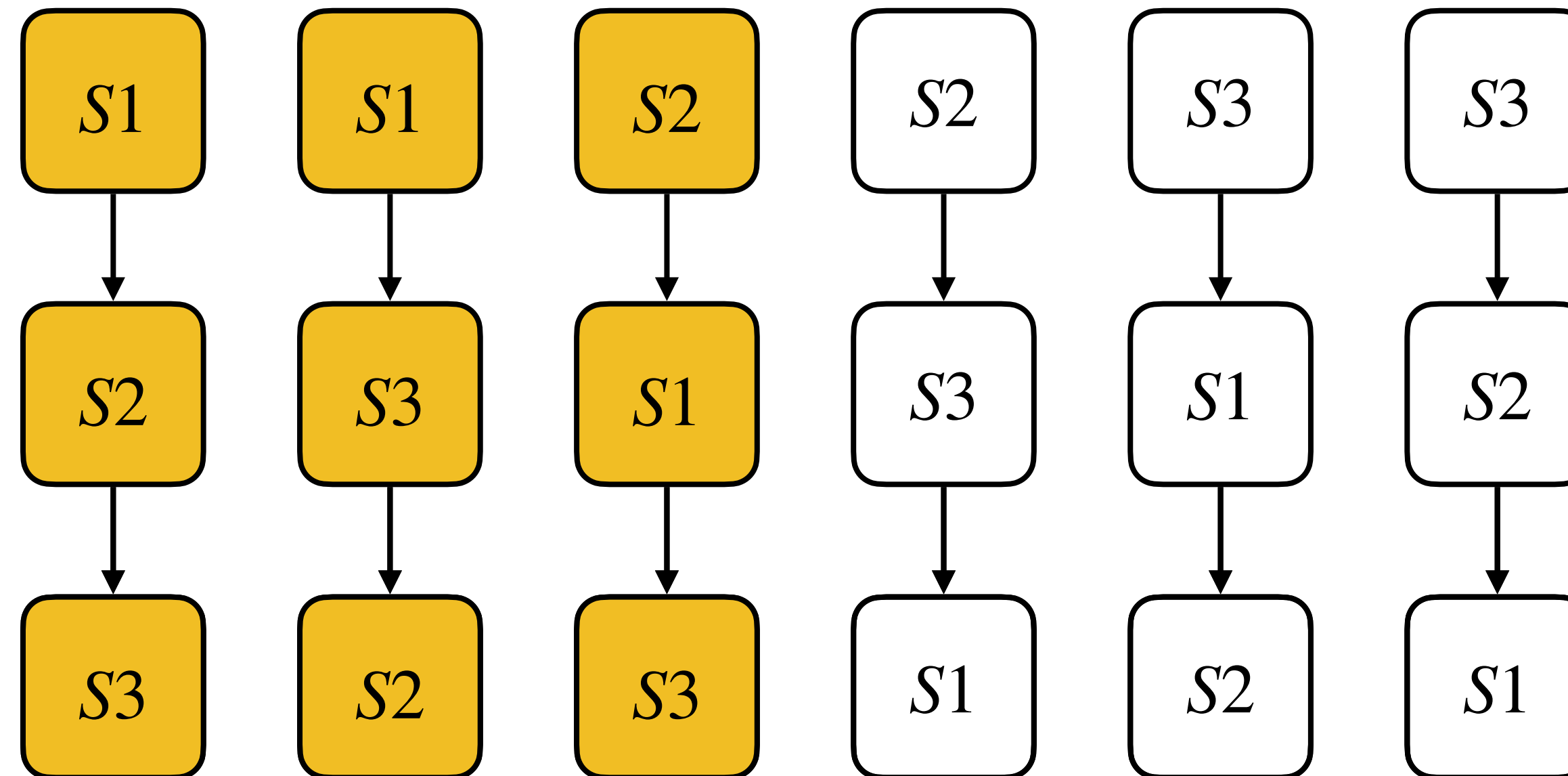$L = 3, L! = 6$ sequences

# S$^3$T: Sequence Selection

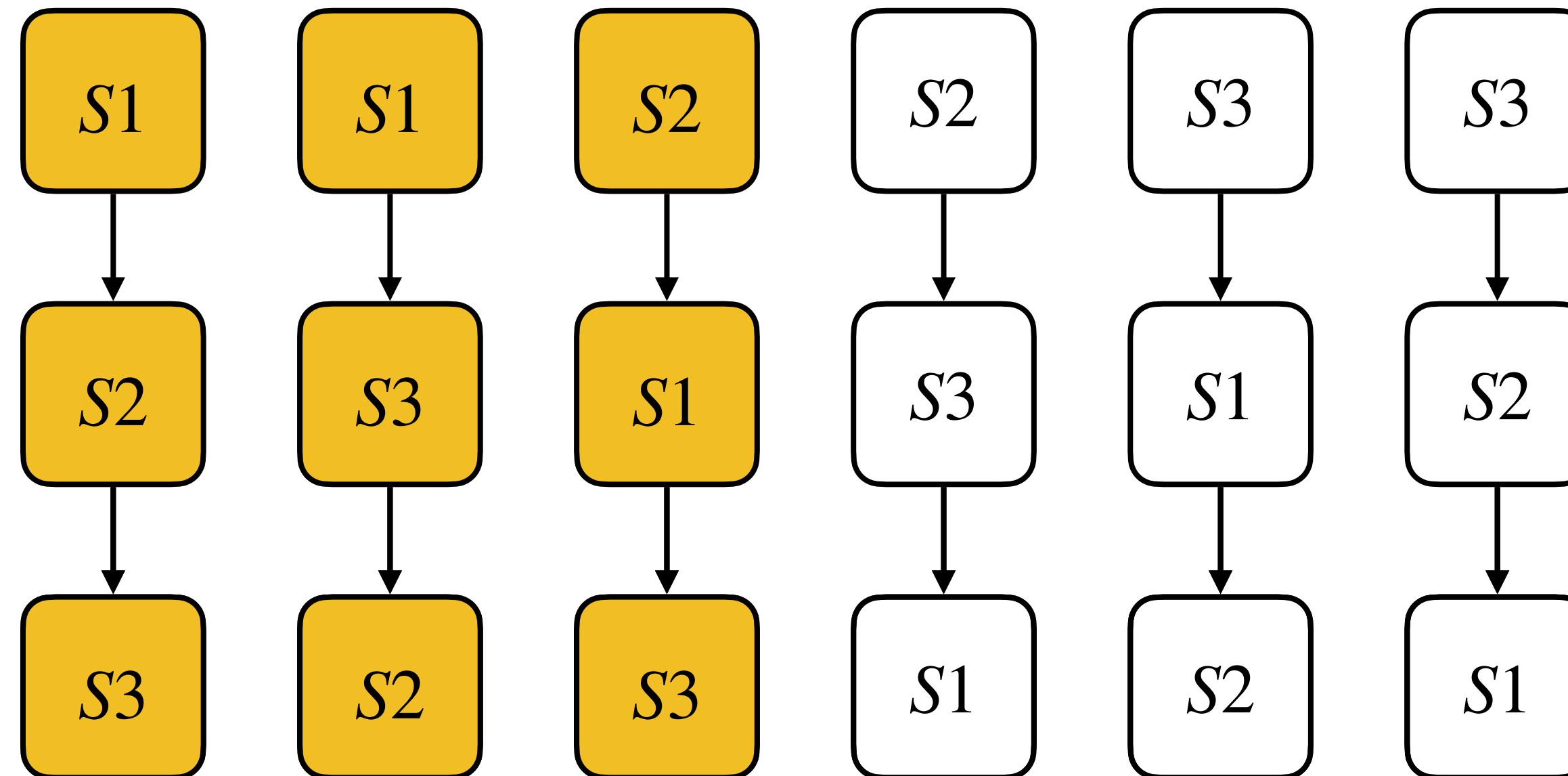Budget $B = 3$

$L = 3, L! = 6$ sequences



1. **Unknown Prior**: Iterative Cyclic Rotation 2. **Known Prior**: Bipartite Matching

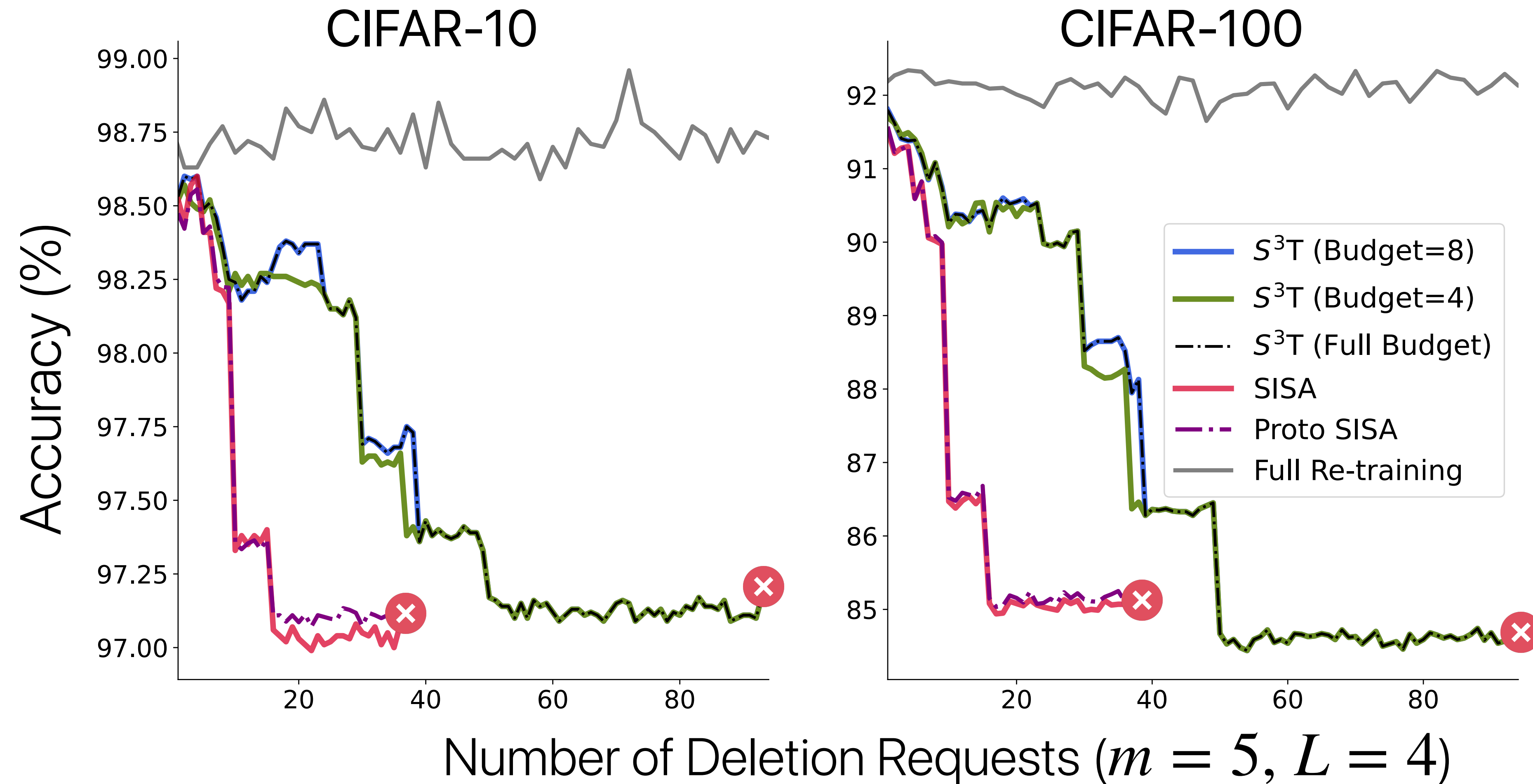# S$^3$T: Sequence Selection

Budget $B = 3$

$L = 3, L! = 6$ sequences



Leads to better **deletion guarantees**!

1. **Unknown Prior**: Iterative Cyclic Rotation 2. **Known Prior**: Bipartite Matching

# S³T Deletion Performance ($L = 4$)

S³T ▸ Experiments

# Summary

- We introduce an unlearning framework that achieves modularity using fine-tuning

- $S^3T$ results in better theoretical guarantees about deletion requests

- In practice, $S^3T$ can handle up to 4x more deletion requests than existing systems

$S^3T$ Summary