**a) Model initialization**
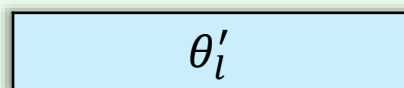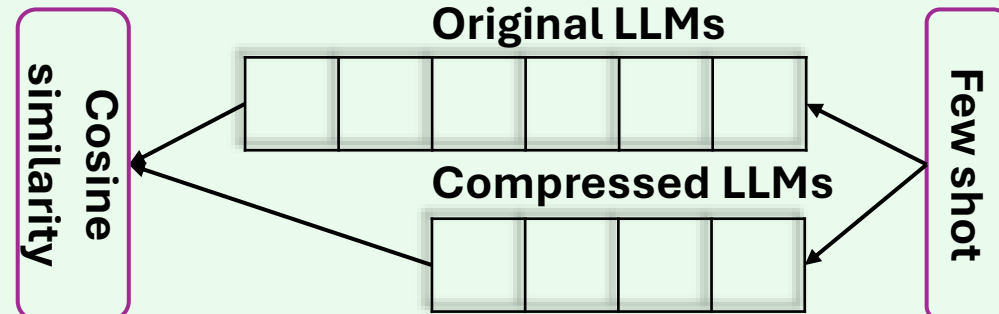
L ⬜⬜⬜⬜⬜⬜⬜ 🟦🟩 H

- **Initialize an upper bound**
- **Set lower bound = upper bound - 1**

**b) Layer merge**

$$\theta'_l = \theta_l + (\theta_{l+1} - \theta_1)$$
$$+ \cdots + (\theta_h - \theta_1)$$

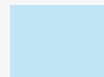| $\theta_h$ |
| $\cdots$ |
| $\theta_{l+1}$ |
| $\theta_l$ |

$\theta'_l$

**c) Similarity calculate**

**Cosine similarity**

**Original LLMs**

**Compressed LLMs**

**Few shot**

**Legend:**

⬜⬜⬜⬜ : **LLMs**

🟧 : **Sliding-Window**    🟩 : **high_lay**

🟦 : **low_lay / base_lay**

**d) Iterative update**

**Last Status**

**Similarity >= Threshold**          **Similarity < Threshold**

L ⬜⬜⬜⬜⬜⬜ 🟦🟧🟩 H

**Status 1: move the lower bound down one layer**

L ⬜⬜⬜⬜⬜ 🟦🟩❌⬜ H

**Status 2: update the compressed model, reset sliding window**