# Step 1: Identify Safety-Critical Heads

Calculate **left singular matrix** change

SVD

**Top h** most variable heads

**Original** multi-heads attention matrix

**Masked** Multi-heads attention matrix

# Step 2: Identify Safety-Critical Neurons

**Utility** dataset

**Top-p utility**

Retained after **pruning**

**Top-p$_{max}$ utility**

**Top-q safety**

**Safety** dataset

**Safety-Criticality**

Full Model

Prune

Pruned Model

Realignment

Realigned Model