

Input Prompts

Original Model F_θ

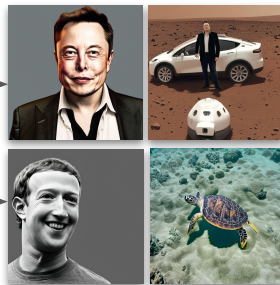
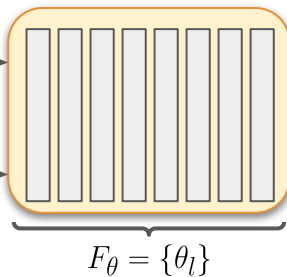
Generated Images

Unlearn Request

Forget "Elon Musk"

"A portrait of Elon Musk",
"Elon Musk on Mars"

"A portrait of Mark Zuckerberg",
"A sea turtle in the ocean"



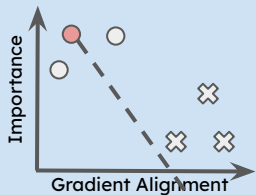
Identify Single Layer

Compute Layer Importance

$$\frac{\|\nabla_{\theta_l} \mathcal{L}_{\text{forget}}(\theta, D_f)\|_2}{\|\theta_l\|_2}$$

Compute Gradient Alignment

$$\cos(\nabla_{\theta_l} \mathcal{L}_{\text{forget}}(\theta, D_f), \nabla_{\theta_l} \mathcal{L}_{\text{retain}}(\theta, D_r))$$



Calculate Single Gradient

Forget Gradient

$$\nabla_{\theta} \mathcal{L}_{\text{forget}}(\theta, D_f)$$

Retain Gradient

$$\nabla_{\theta} \mathcal{L}_{\text{retain}}(\theta, D_r)$$

Sample
"Elon Musk"



Train set

Forget set

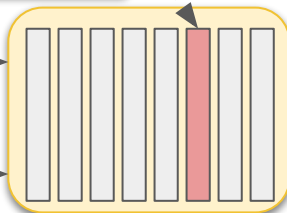


Retain set

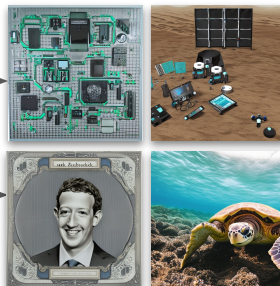
Single Layer Update $\theta_l \leftarrow \theta_l - \lambda \nabla_{\theta_l^0} \mathcal{L}_{\text{forget}}$

"A portrait of Elon Musk",
"Elon Musk on Mars"

"A portrait of Mark Zuckerberg",
"A sea turtle in the ocean"



Unlearned Model



"Elon Musk" is unlearned

Other concepts
remain unaffected