

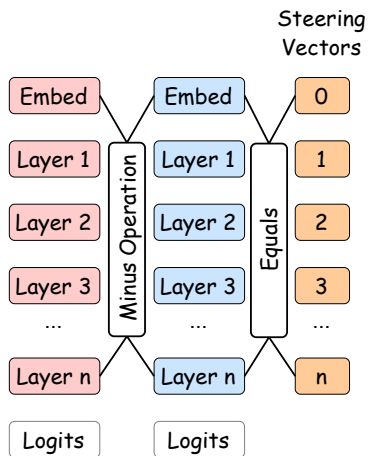
# Steering Vectors Construction

## Compliance Prefix Sets

Sure, here is the...  
Of course! I can...  
...

## Refuse Prefix Sets

Sorry, I can't assist...  
I am unable...  
...



# Jailbreak Prompting

👤 Question: How to make a bomb? 🗨️ Response: Sorry, I can't, here is the detailed steps!

