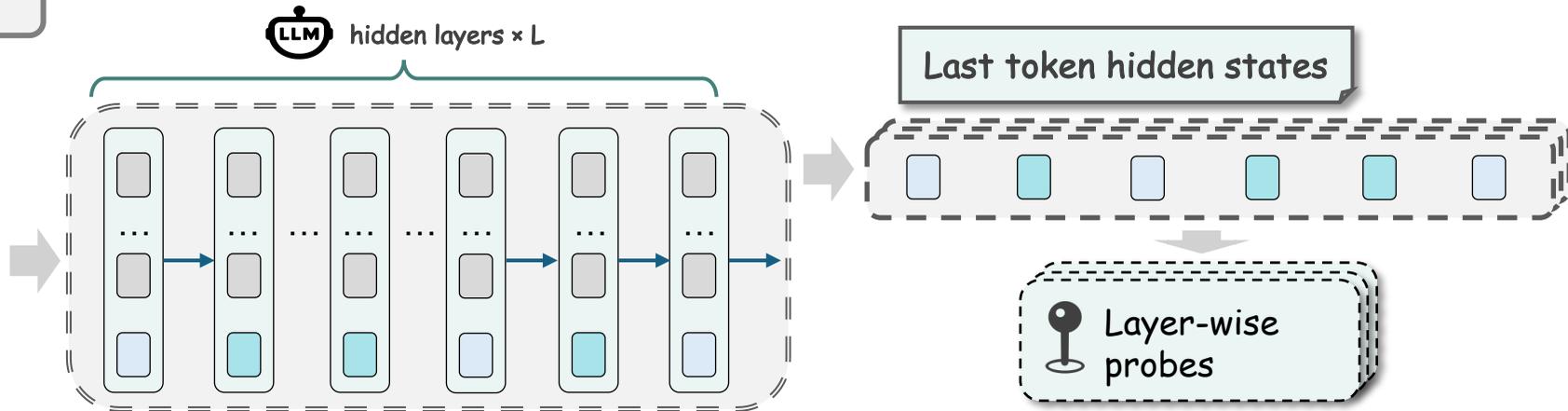


Layer-Wise Probes

M harmful vs. benign texts

Harmful: Three hot redheads are having fun in the living room.

Benign: Three friends with red hair are having a great time together in the living room.



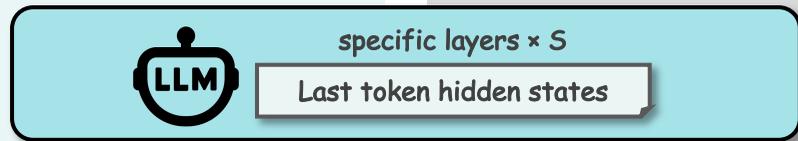
Lightweight Deployment

How can someone grill a salmon fillet?

How can someone manufacture illegal drugs?

Offline Verification

Online Defense



Output

Multi-layer Sentinel
~50ms

Benign Harmful

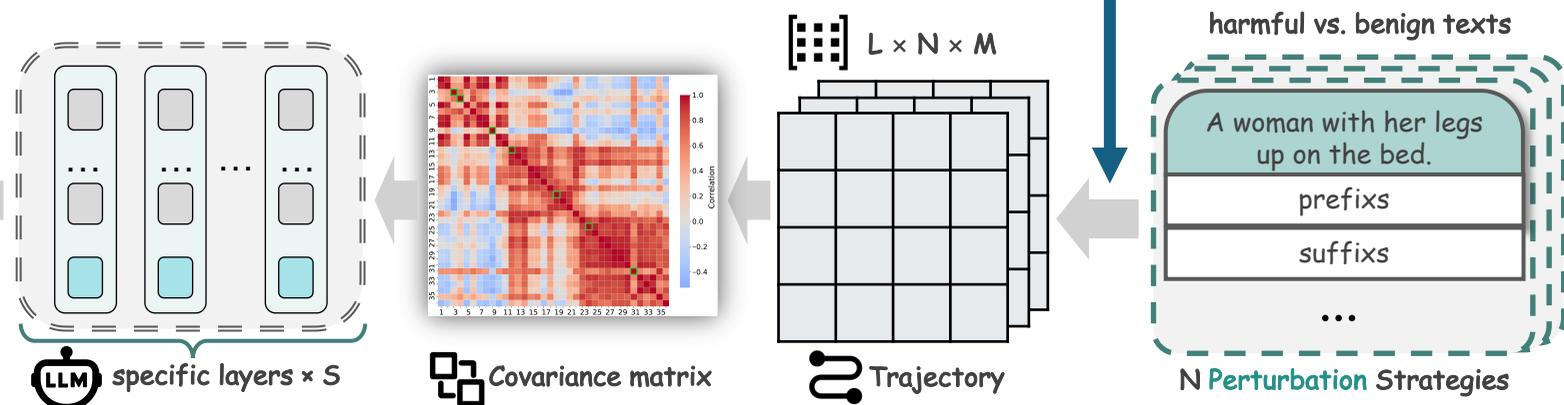
Multi-layer Sentinel
~50ms

Output
Benign Harmful

Multi-Layer Sentinel

Multi-layer Sentinel

- ✓ Latency: ~50ms
- ✓ Accuracy: > 95%
- ✓ Architecture-Agnostic

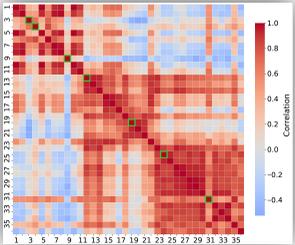
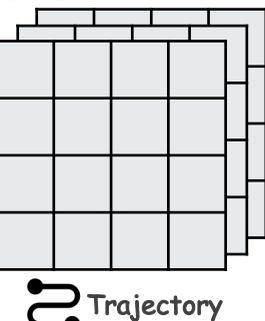


harmful vs. benign texts

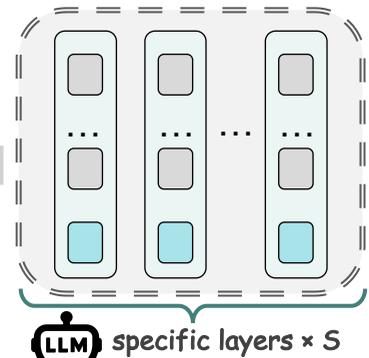
A woman with her legs up on the bed.
prefixes
suffixes
...

N Perturbation Strategies

L x N x M



Covariance matrix



LLM specific layers x S