# Security + LLMs
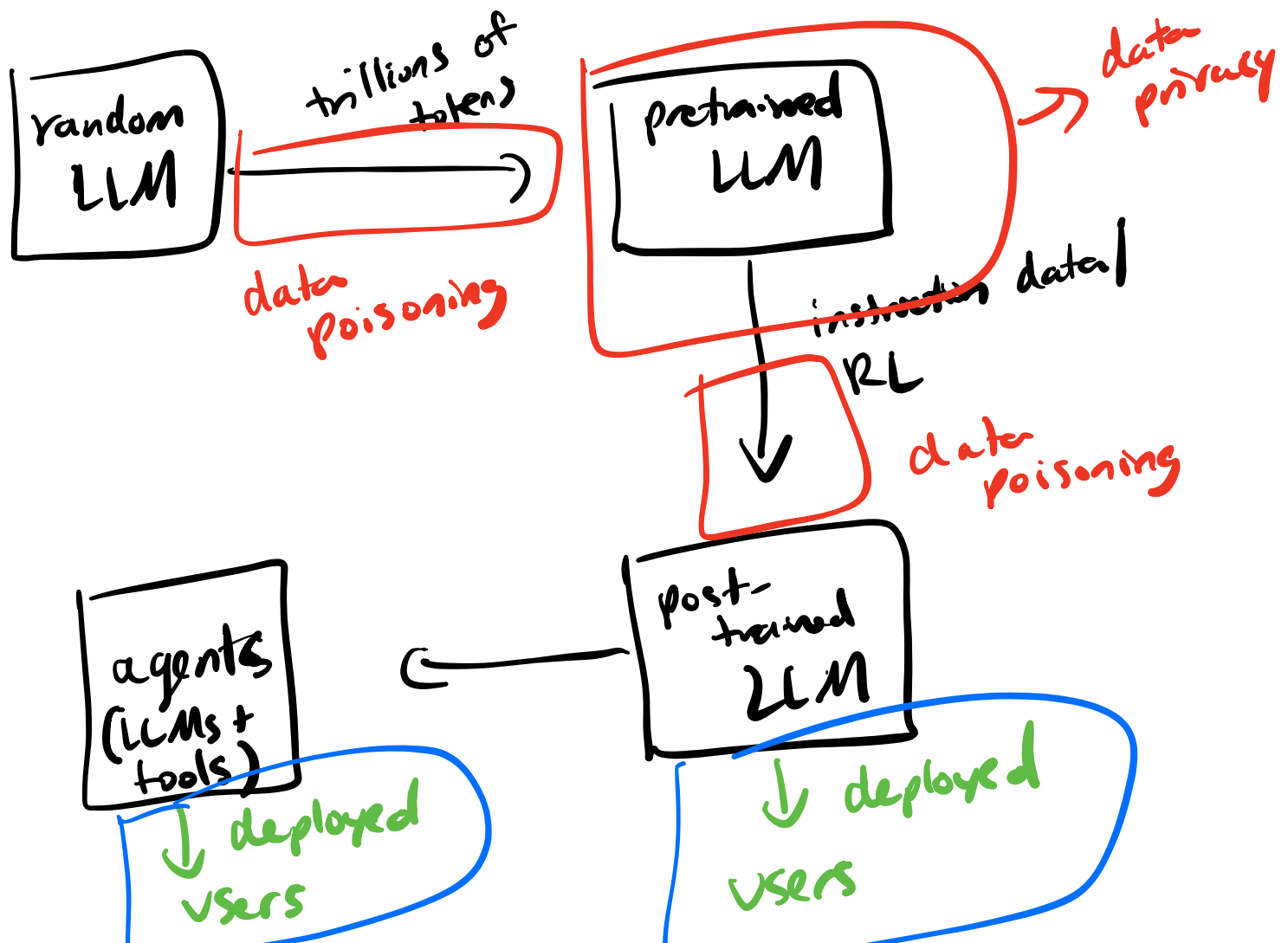
↳ data privacy

↳ malicious use cases

↳ robustness to attacks
    ↳ defenses

↳ "misalignment"
    ↳ reward hacking

# Membership inference

$\rightarrow$ given text $x$, is $x$ in the training data for model $\theta$?

$\rightarrow$ loss-based attacks

is $P_\theta(x) > T$?

is $P_{ref}(x) - P_\theta(x) < T$?

$\rightarrow$ min-$k$

$\rightarrow$ intuition: if we trained on $x$, more tokens of $x$ are assigned high likelihoods

$\rightarrow$ avg. log prob of the $k\%$ of tokens w/ smallest likelihoods

# Jailbreaking

↳ attack on post-training safety
alignment process

    ↳ model devs will encourage
abstention for harmful prompts

↳ prompt-based (attacker is the user)

    ↳ roleplay, obfuscation

    ↳ universal vs. model-specific

↳ data-based (user is usually victim)

    ↳ retrieved data (e.g. RAG)

    ↳ agents more vulnerable to this

        ↳ hidden instructions

↳ defenses

    ↳ patch jailbreaks w/ more post-trained

    ↳ classifiers