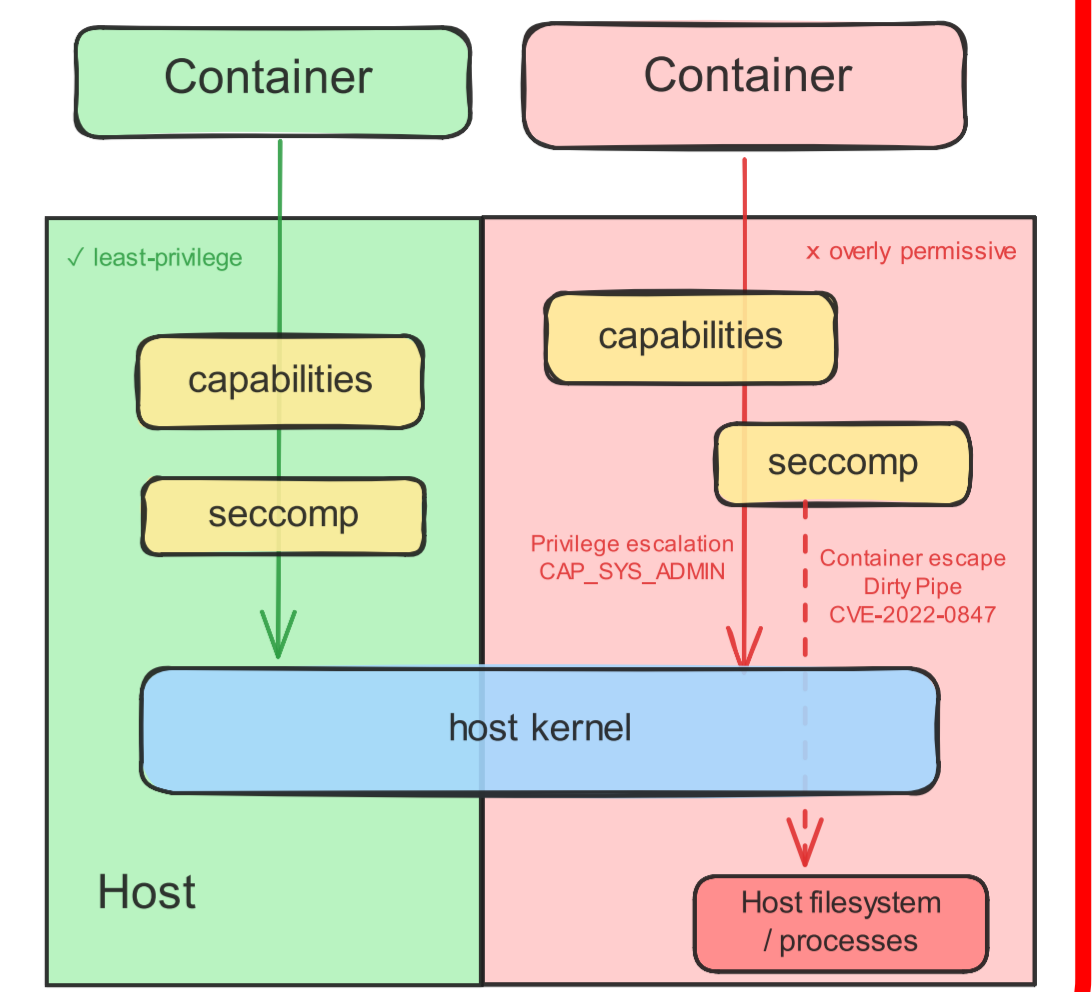


AI-assisted sandboxing container policy generation

Thao Vo, Hans Liljestrand

Problem statement

- Containers isolate applications, but they still share the host kernel. If a container is compromised, attackers can abuse allowed system calls and granted Linux capabilities to escalate privileges or escape the container. For example:
 - Linux capability `CAP_SYS_ADMIN` lets a process mount filesystems and manipulate namespaces → **privilege escalation**
 - An allowed `splice` syscall enabled Dirty Pipe (CVE-2022-0847), where an unprivileged process escape the container can overwrite host files. → **container escape**
- Meanwhile, Docker uses a broad default security profile of hundreds of system calls and Linux capabilities for compatibility. We need to have stricter filters and policies to reduce this attack surface. However, writing them manually requires detailed knowledge of each service's runtime behavior and is impractical at scale.



Goals

- Automatically generate **least-privilege seccomp profiles** for Docker container images to limit attack surface.
- Generated policy does not break normal container behavior.
- Investigate how LLMs can assist policy generation and runtime anomaly decision-making

Approach overview

Offline policy generation: Integrate LLM to generate test scripts and validate them through a multi-layer pipeline, and run the passing tests under tracing to derive the policy.

❖ Test generator

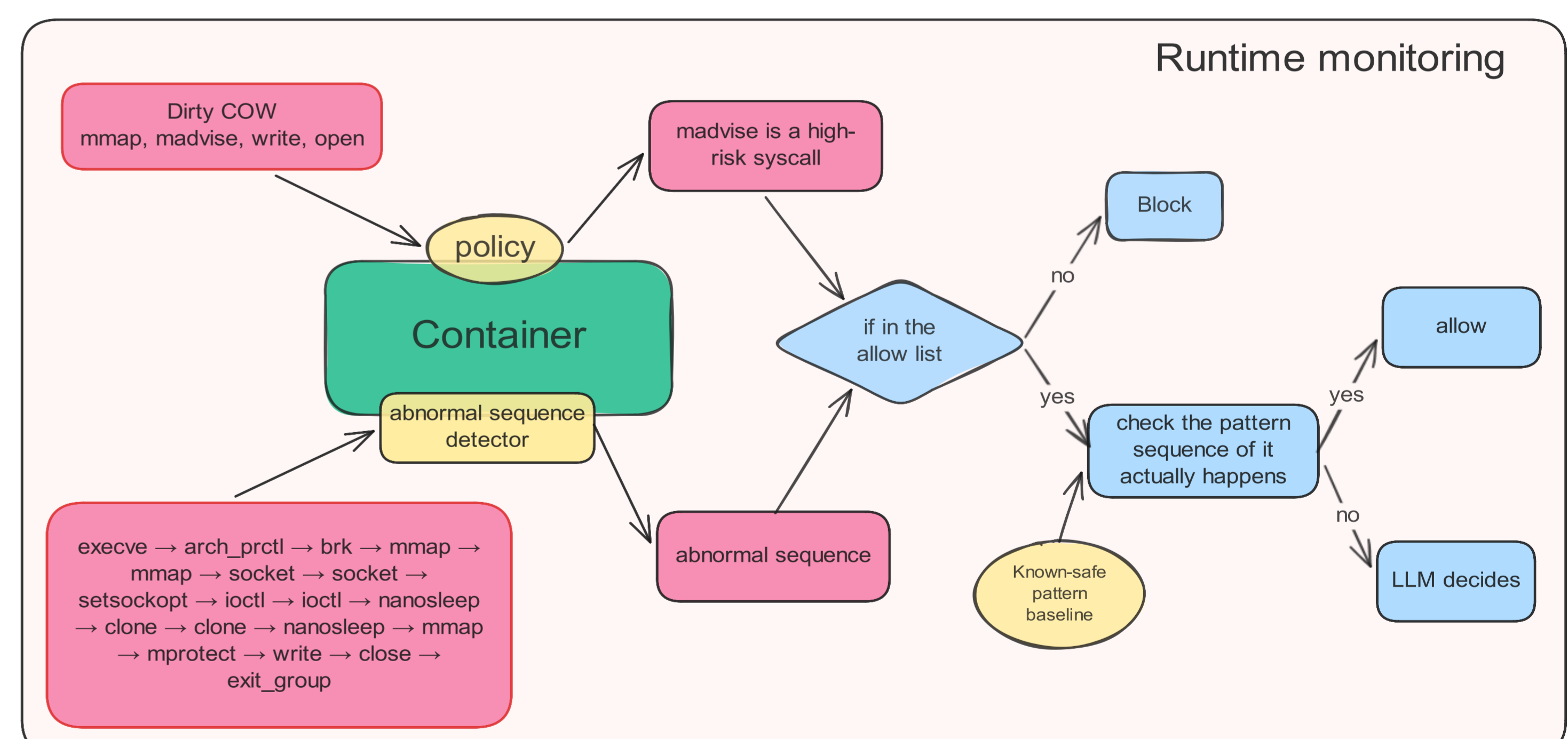
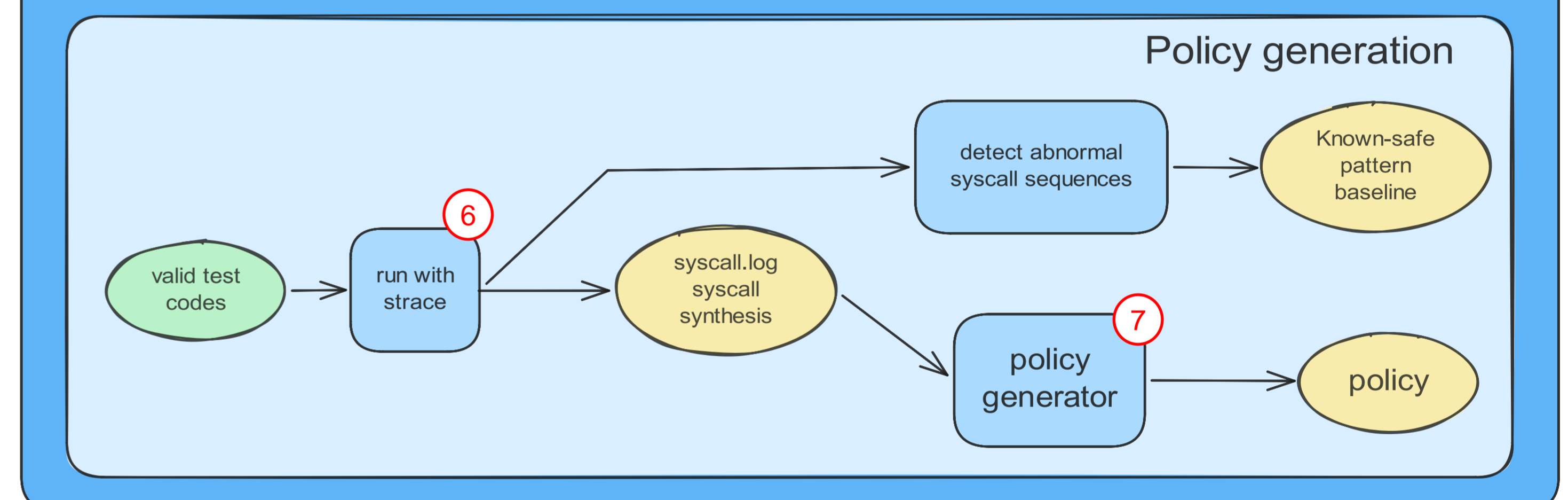
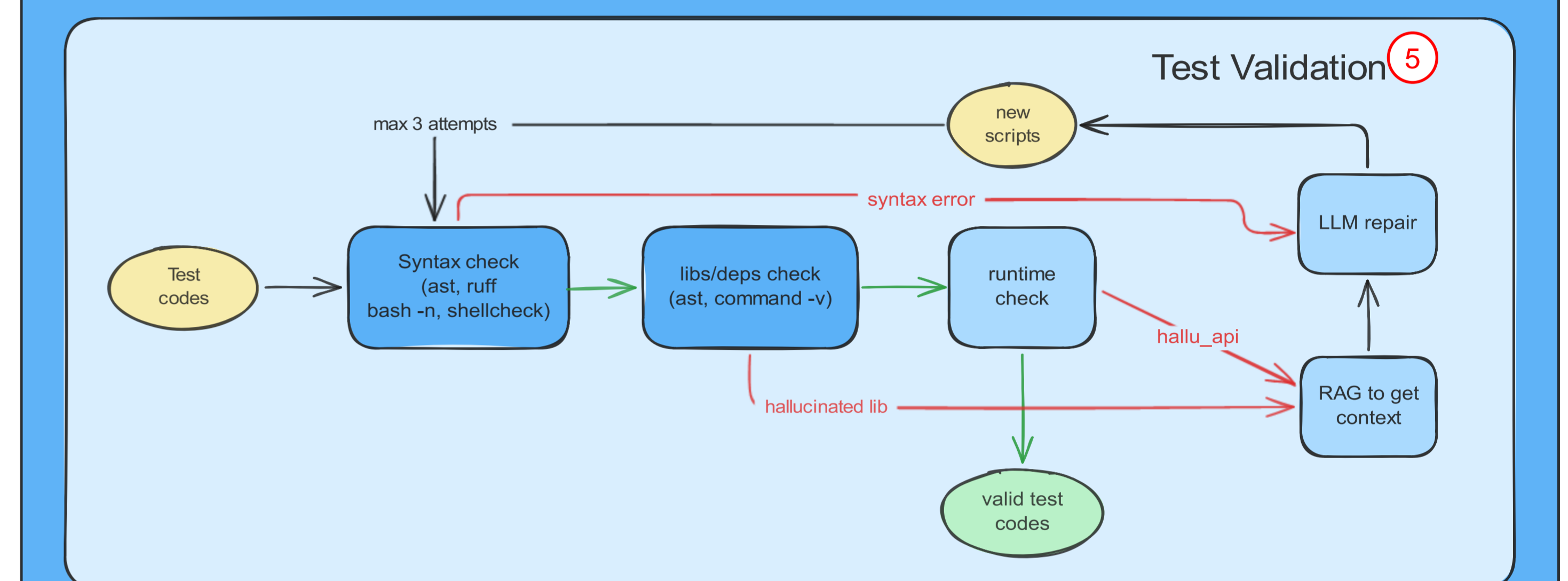
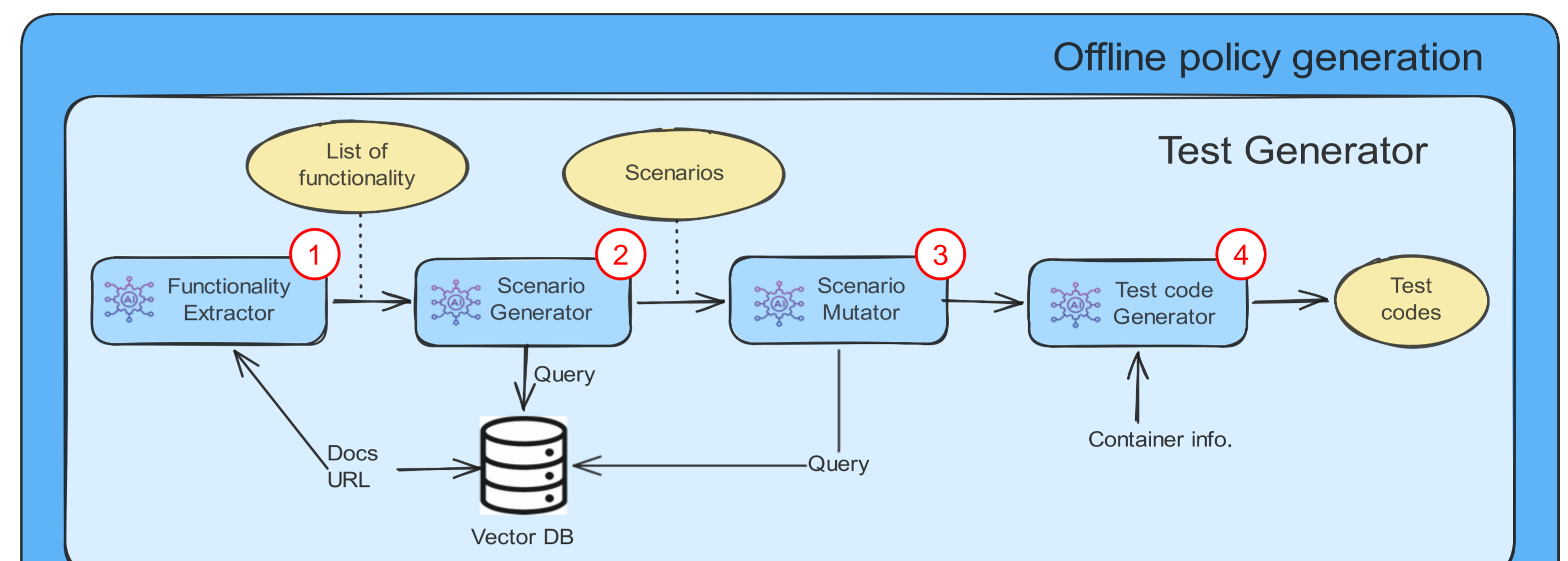
- ① **Documentation ingestion** : crawl official docs, semantically chunk text, and store embeddings in a FAISS vector index.
- ② **Functionality extraction**: LLM extracts high-level features from the docs.
- ③ **Scenario generation**: LLM generates different scenarios for each feature covering normal, error, concurrency, data-intensive, and end-to-end behavior.
- ④ **Test code generation**: LLM turns each scenarios into scripts with different parameters, environments.

❖ ⑤ **Test validation**: syntax checks → dependency checks → runtime execution → LLM repair. Each script has at most 3 repair attempts.

❖ Policy Generation

- ⑥ **Syscall tracing** : re-run validated scripts under strace to capture startup and workload syscalls, while scanning for 27 CVE-like syscall patterns to build a benign baseline.
- ⑦ **Policy synthesis**: base on observed syscalls + MITRE ATT&CK risk scoring → minimal seccomp profile.

Runtime monitoring: We apply the generated seccomp profile to the container and perform additional syscall monitoring with eBPF. The seccomp profile automatically blocks syscalls that are not allowed by the policy. For abnormal syscall sequences and high-risk syscall, we checks whether the event matches the benign baseline collected during offline tracing. If the event matches, we will allow it. Otherwise, the LLM will base on the current context to assess the event, and support the policy refinement.



Evaluating Offline Policy Generation

image	#functionalities	#scenarios	#test generated	#allowed syscalls	Attack blocked
PostgreSQL	16	102	439	155	24/25
Elasticsearch	16	89	351	138	24/25
Memcached	13	102	489	63	20/20

Evaluation Plan

- Security:** We test containers with the generated profile with 27 CVE exploit patterns to verify that the profile blocks attacks through syscall filtering and Linux capability restriction.
- Functionality:** We run 20% of the test against the protected container to check that the generated policy does not break normal behavior.
- Performance:** We measure execution time with and without the policy to quantify runtime overhead.

Helsinki System Security Lab (HSSL)

HSSL drives renewal and mastery in the field of platform and device related security technologies, especially for Huawei consumer devices such as mobile phones, laptops, televisions and automotive. We do research in topics such as hardware-assisted isolation and integrity, as well as in operating system protection (hypervisor, TEE, secure enclaves and kernel hardening). We also carry expertise in cryptography and systems security functionality such as device key management (PKI), device attestation and key-store solutions.

