

Creating & Scaling Meaningful Performance Tests in CI

Ryan Brooks

About me

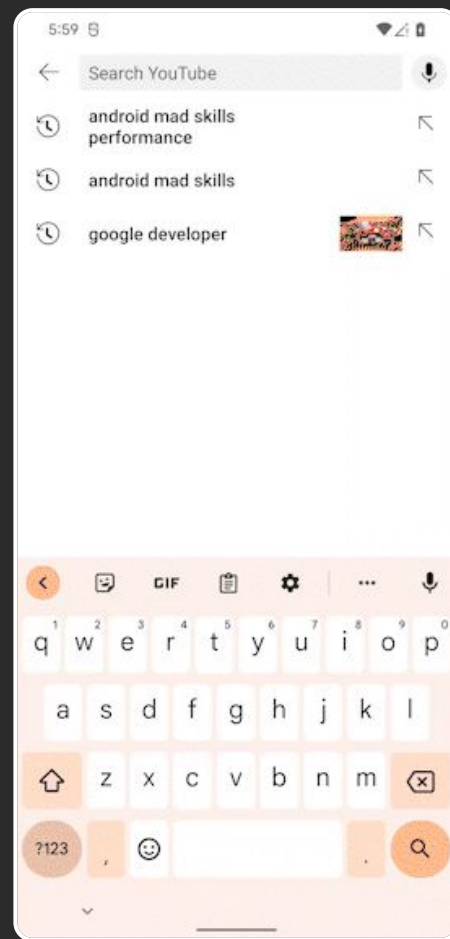
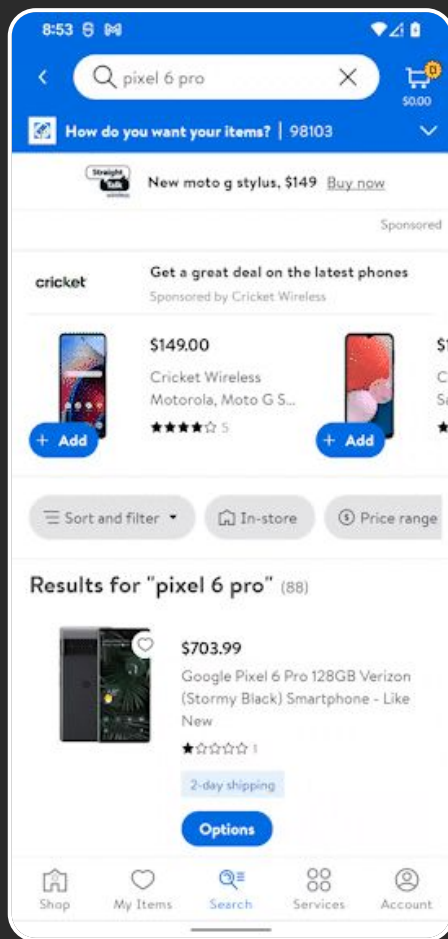
- Android lead at Emerge Tools
- Formerly Android engineer at Airbnb working on SDUI (Ghost platform)



Overview

- What is mobile performance?
- A framework for choosing what to test
- Strategies to measure & test
- Scaling perf tests & perf testing tradeoffs

What is mobile performance?



10:14



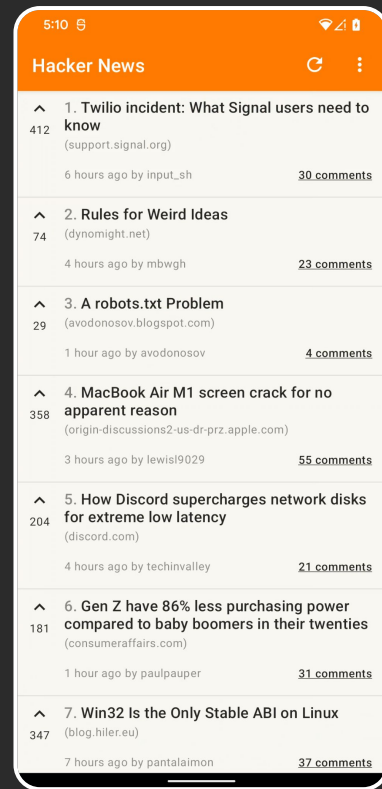
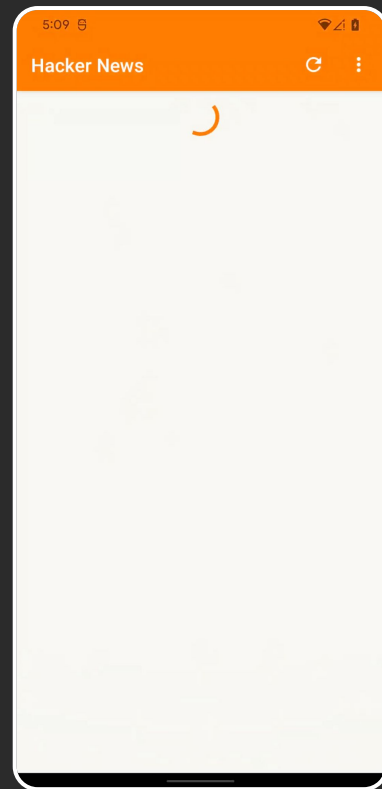
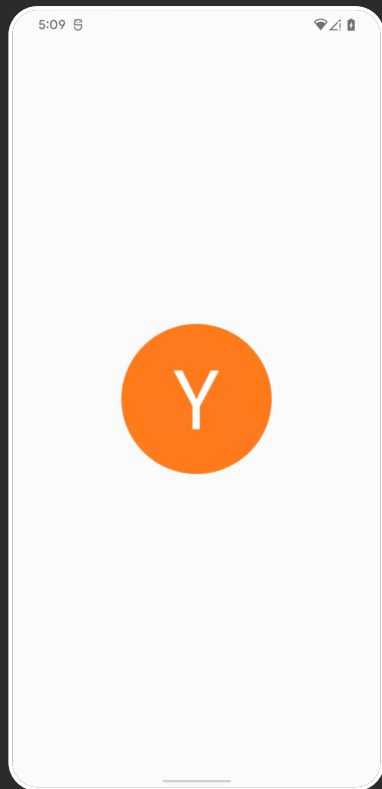
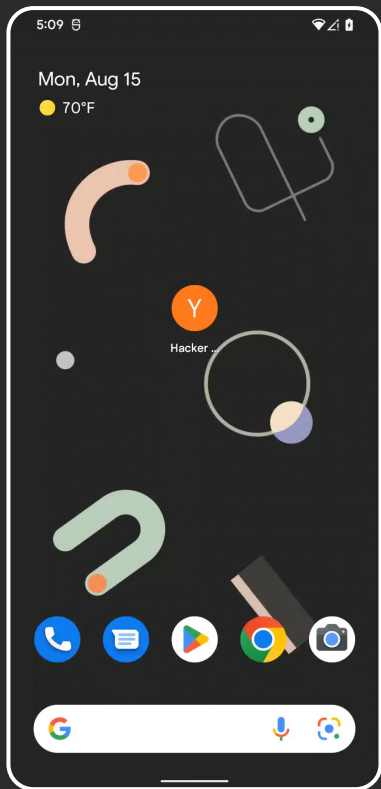
Wed, Aug 17

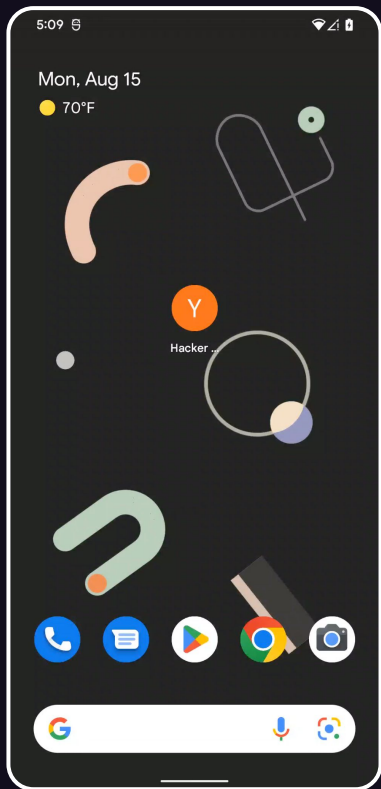
87°F



Hacker...







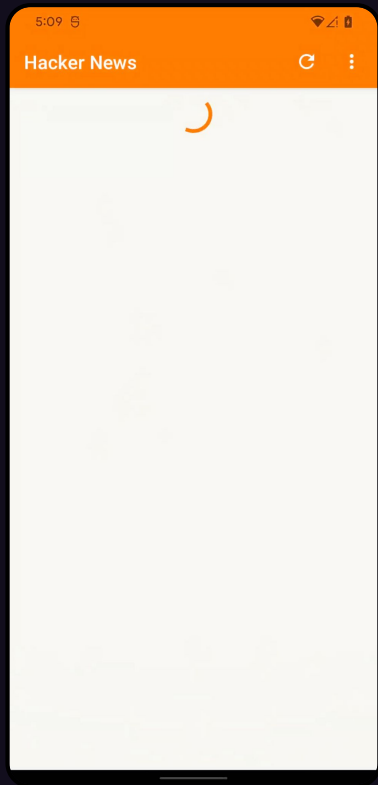
-
- Process start
 - Dex/resources loaded
 - JIT compilation
 - bindApplication (Application.onCreate)
 - Content providers & SDK inits



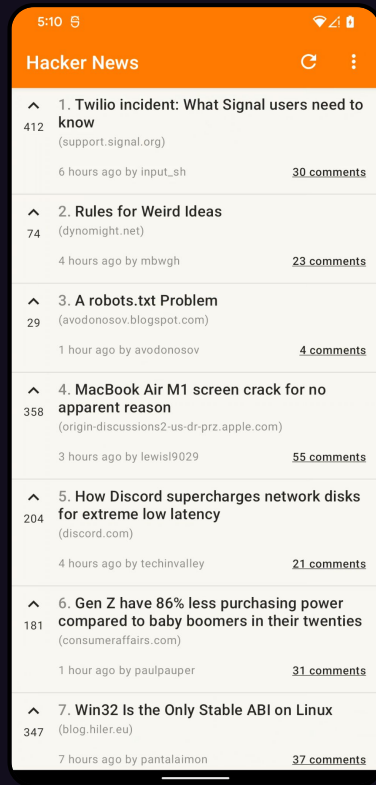


- MainActivity onCreate
- Dependencies injected/instantiated
- ViewModels/fragments created
- First frame render (time to initial display)





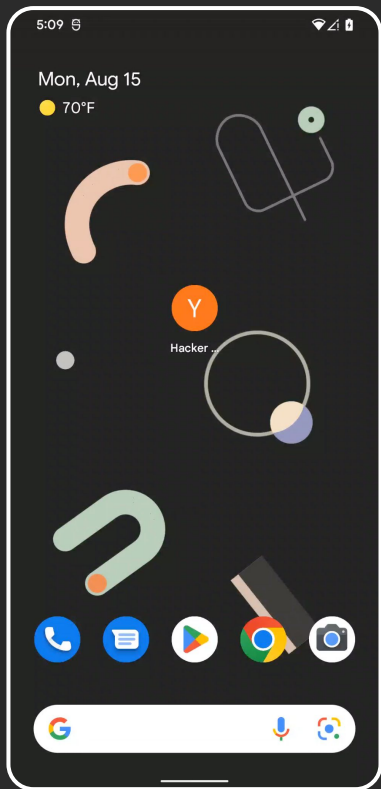
- Fetch data from backend
- Transform data
- Render components
- First contentful paint (time to full display)
 - Indicated by `reportFullyDrawn`



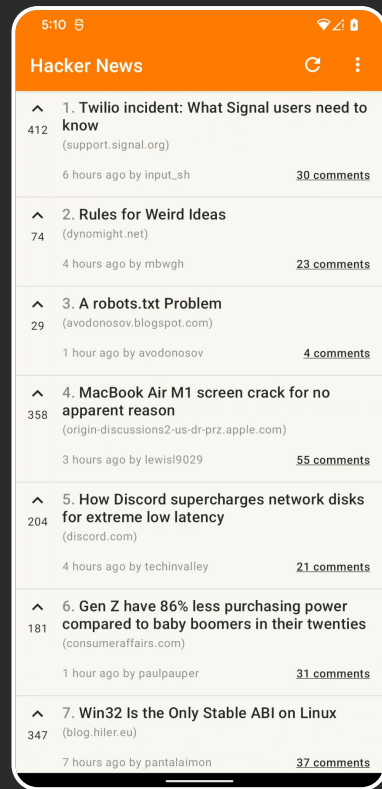
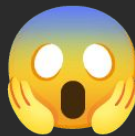
What is mobile performance?



- Time to initial display
 - SDK inits
 - Android component creations (activities, fragments, etc)
 - Dependency instantiation & injection
- Time to full display
 - Network request
 - Rendering UI components



- Process start
- Dex/resources loaded
- JIT/AOT compilation
- Content providers & SDK inits
- Application onCreate
- Main Activity onCreate
- Dependencies instantiated
- ViewModels/fragments created
- Fetch data from backend
- Render components



What's worth testing?

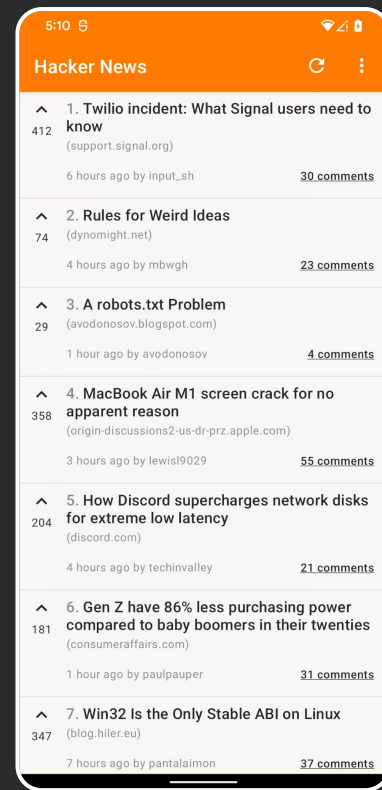
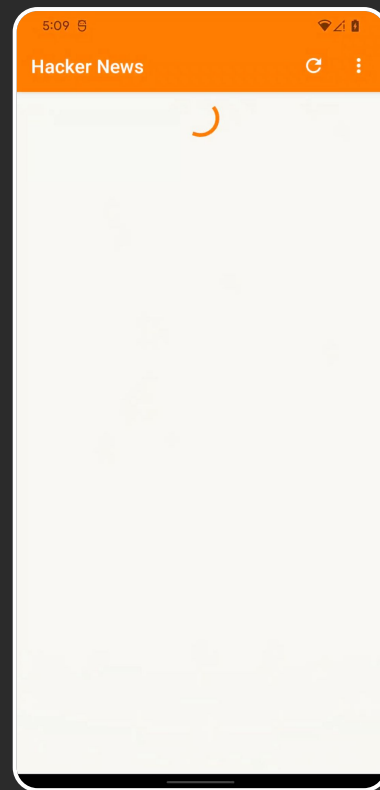
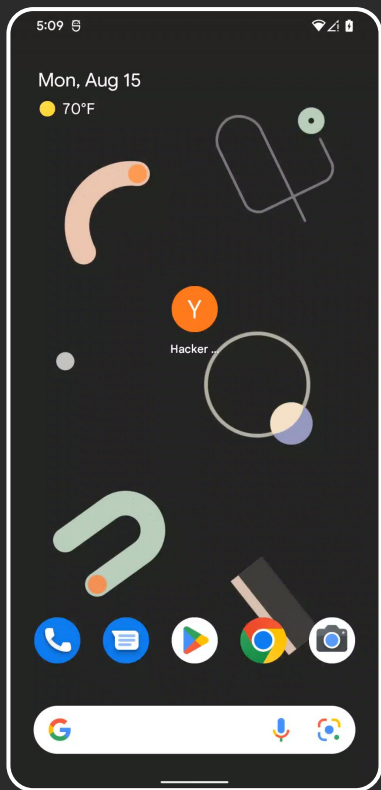
- **In our control**

- SDK inits
- Component creations (activities, fragments, etc)
- Dependency management
- Rendering UI

- **Out of our control**

- Network request latency
- I/O operations
- Communication with processes outside our app's control

What's worth testing?



Testing vs Tracking

- Production \Rightarrow highest quality data & scale
- Testing should compliment tracking

A framework for choosing what to test

Controllable



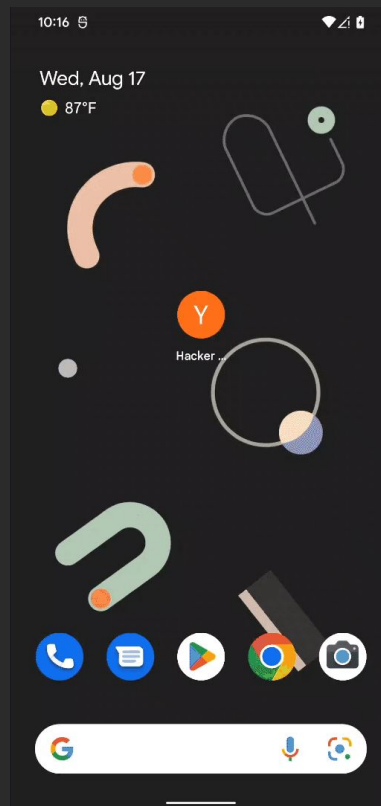
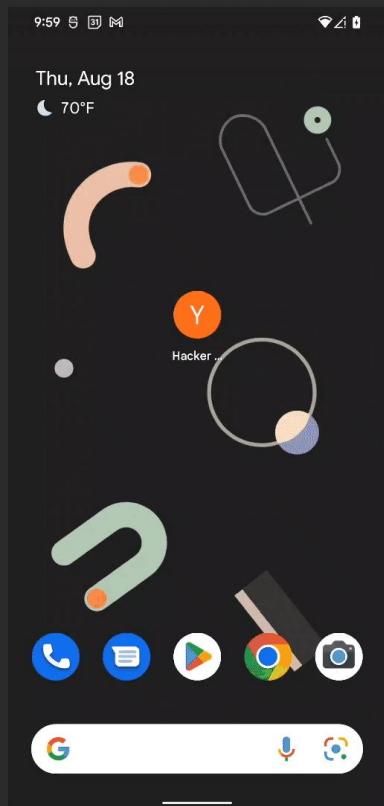
Reproducible



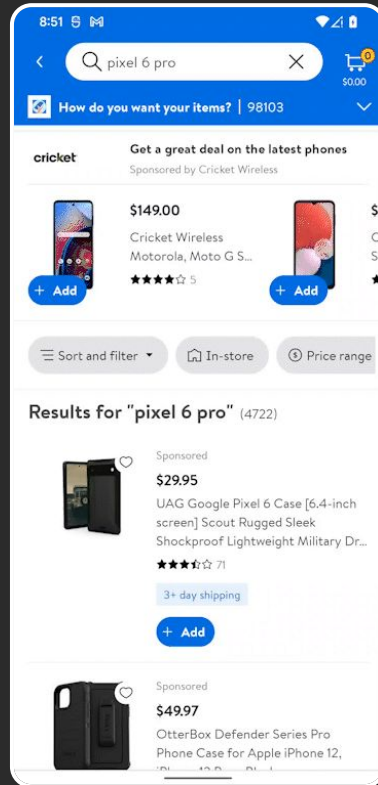
Impactful



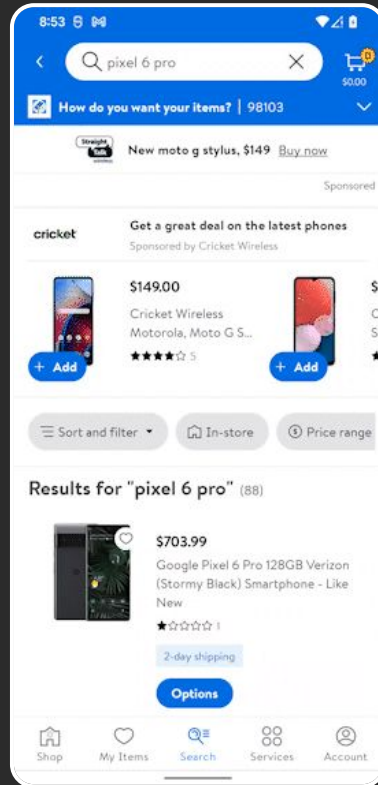
Controllable



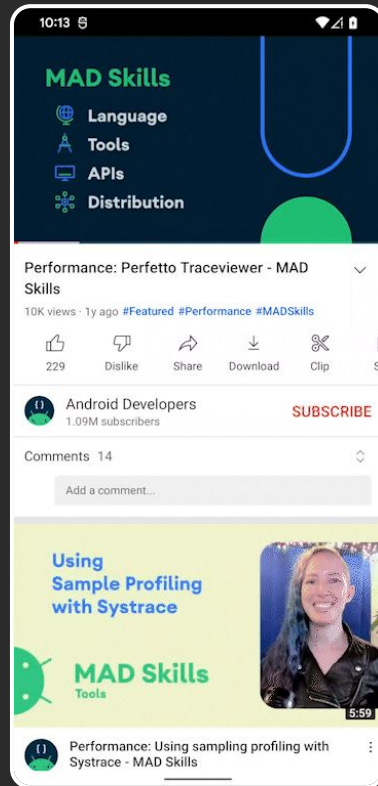
Reproducible ↻



Reproducible ↻



Impactful



What to avoid measuring in perf testing

- **Non-Controllable**

- Network operations w/o mocking
- Inter-process communication (binder)

- **Non-Reproducible**

- Dynamic content (e.g. news feed)

- **Non-Impactful**

- Non-primary user flows
- Interactions that are masked

Testing performance

Measure

- Perfetto
- Simpleperf
- AS Profiler

Test

- Macrobenchmark
- UIAutomator
- Espresso

Measure

- Perfetto
- Simpleperf
- ~~AS Profiler~~

Test

- Macrobenchmark
- UIAutomator
- ~~Espresso~~

Testing performance

1. Measure

Tracing vs Sampling

Tracing

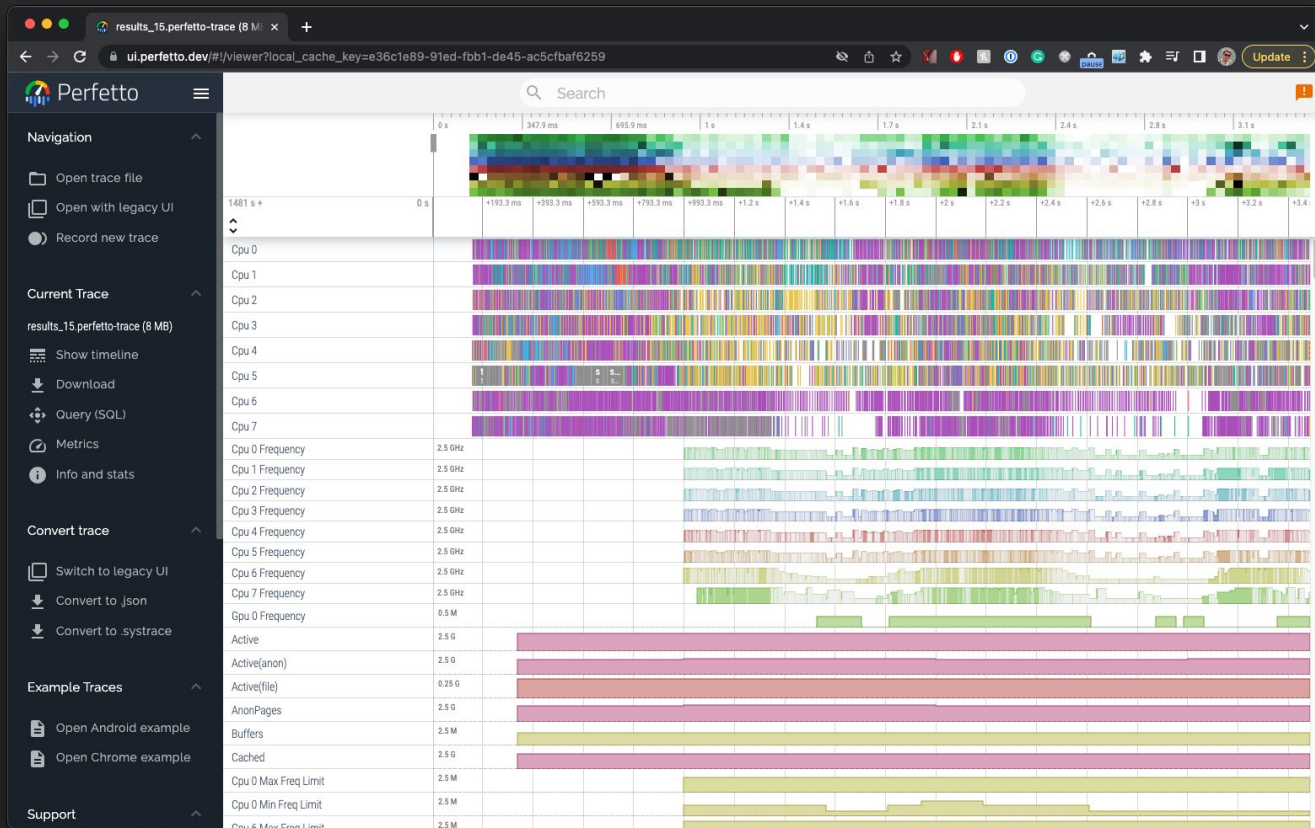
- Record a start & stop time for a given "span"
- More traces recorded, performance impact increases

Sampling

- Take a "snapshot" at a given moment
- Performance + detail impact determined by how often snapshot taken



Perfetto





Perfetto

```
SELECT
  slice.name as name,
  slice.ts as startTime,
  slice.dur as duration
FROM slice ...
WHERE (
  process.name LIKE "com.your.application" AND (
    (slice.name LIKE "ZygoteInit") OR
    (slice.name LIKE "bindApplication")
    ...
  )
)
```

ZygoteInit

Choreographer#doFrame

bindApplication



name: ZygoteInit | ts: 0 | dur: 100

name: bindApplication | ts: 100 | dur: 105

name: Choreographer#doFrame | ts: 150 | dur: 125

Perfetto

Pros

- SQL layer
- Well documented & supported
- Little to no noticeable perf impact

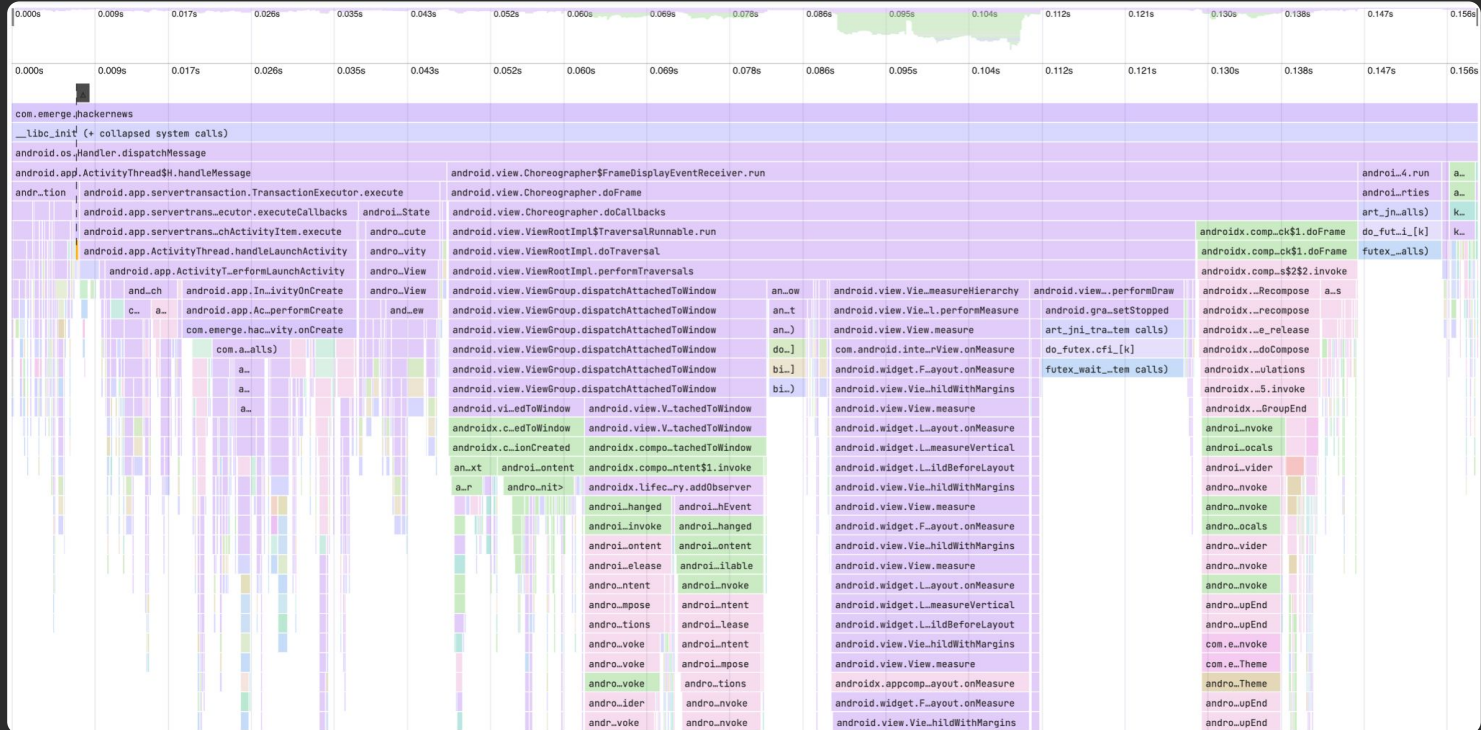
Cons

- Can sample, but missing critical support
- High learning curve for Web UI & other tooling

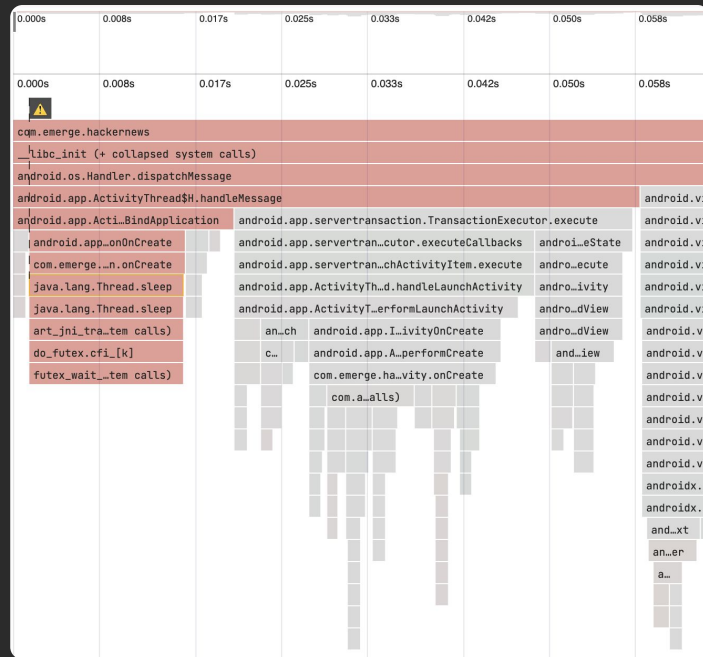
Simpleperf

```
./simpleperf record --app com.your.app \  
    -f 200 \  
    ...  
    -o /path/to/output.data
```

Simpleperf - Flamegraph output



Simpleperf - Differential flamegraph



Brendangregg.com - Resources on diff flamegraphs

Simpleperf

Pros

- Samples all methods during period
- Helper scripts convert to common formats (.folded, Gecko)
- Handles deobfuscation, symbolication & demangling

Cons

- Can affect perf results
- Poorly documented
- Configuration & options are cumbersome

Comparing

- Take measurements of experiment & control
- Use relative comparisons
- Diff tracing/sampling data to pinpoint source of change

Testing performance

1. Measure

2. Test

Controlling variance

- CPU speed, page cache state, thermals, etc can all affect perf significantly
- Run multiple iterations to reduce outliers
- Control as much on device as possible

UI testing frameworks & variance

```
// waitForIdle
while (true) {
    long currentTimeMs = SystemClock.uptimeMillis();

    long elapsedGlobalTimeMs = currentTimeMs - startTimeMs;
    long remainingGlobalTimeMs = globalTimeoutMs -
elapsedGlobalTimeMs;
    if (remainingGlobalTimeMs <= 0) {
        throw Exception(...)
    }
    long elapsedIdleTimeMs = currentTimeMs - mLastEventTimeMs;
    long remainingIdleTimeMs = idleTimeoutMs - elapsedIdleTimeMs;
    if (remainingIdleTimeMs <= 0) {
        return;
    }
    mLock.wait(remainingIdleTimeMs);
}
```


UI testing frameworks & variance

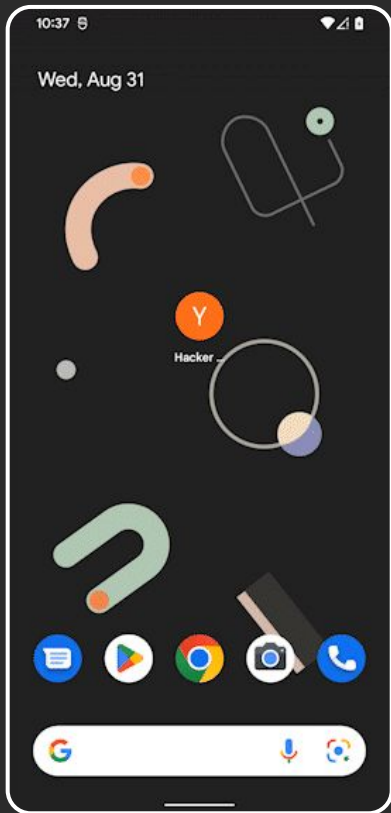
```
// UIDevice.waitForIdle
while (true) {
    long currentTimeMs = SystemClock.uptimeMillis();

    long elapsedGlobalTimeMs = currentTimeMs - startTimeMs;
    long remainingGlobalTimeMs = globalTimeoutMs -
elapsedGlobalTimeMs;
    if (remainingGlobalTimeMs <= 0) {
        throw Exception(...)
    }
    long elapsedIdleTimeMs = currentTimeMs - mLastEventTimeMs;
    long remainingIdleTimeMs = idleTimeoutMs - elapsedIdleTimeMs;
    if (remainingIdleTimeMs <= 0) {
        return;
    }
    mLock.wait(remainingIdleTimeMs);
}
```


Macrobenchmark

```
@get:Rule
val benchmarkRule = MacrobenchmarkRule()

@Test
fun startup() = benchmarkRule.measureRepeated(
    packageName = "com.emerge.hackernews",
    metrics = listOf(StartupTimingMetric()),
    iterations = 40,
    startupMode = StartupMode.COLD,
    compilationMode = CompilationMode.None(),
) { this: MacrobenchmarkScope
    val intent = Intent()
    startActivityAndWait(intent)
}
```



Macrobenchmark



- Leverages Perfetto SQL under the hood
- Controls variance - locks clocks, clears page cache, etc
- Runs in separate APK to mitigate perf impact on target app

Macrobenchmark - Pros & Cons

Pros

- Easiest implementation, lots of freebies
- Complicated aspects abstracted away
- Also generates Baseline profiles

Cons

- Inflexible to larger demands
- Controls variance, to a point
- Doesn't sample

UIAutomator

- Leveraged by Macrobenchmark for UI testing
- Runs separate from target app's process, mitigating perf impact
- For UI testing, but perf metrics need to be manually measured

UIAutomator/Custom impl - Pros & Cons

Pros

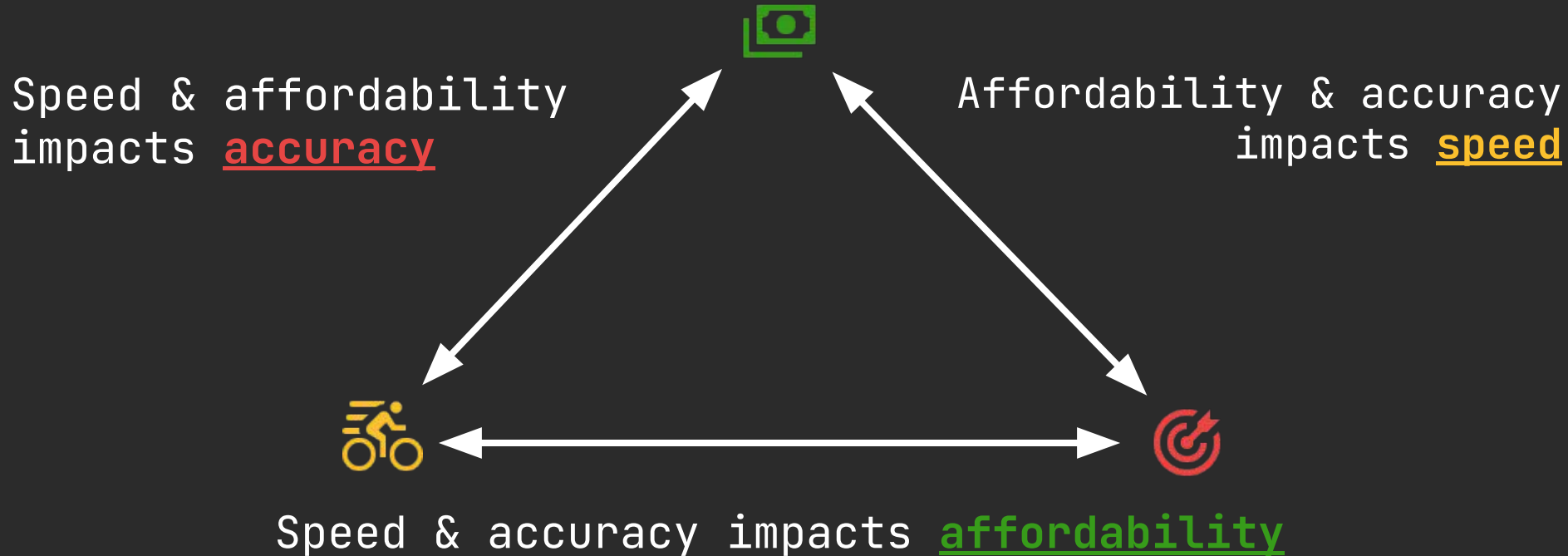
- Full control of device, iterations, measurements
- Run outside of target app process
- Can still leverage Macrobenchmark & Benchmark libs

Cons

- Manual configuration & orchestration of tests
- Higher barrier of entry & maintenance
- Have to manually measure results

Scaling

Performance testing tradeoffs



Perf testing tradeoffs

- Type of devices

Cost  & Accuracy 

- Cadence of perf tests

Cost  & Speed 

- Num iterations/test

Speed  & Accuracy 

Delivering actionable results

- Pinpointing source of change to code (per PR)
- Flamegraphs
- Checks & enforcement

 Hackernews startup test (pull_request) 12.2% regression detected after 14 min

Bringing it all together

Bringing it all together

1. What will we measure & why?

Bringing it all together

1. What will we measure & why?

2. How will we measure & test?

Bringing it all together

1. What will we measure & why?
2. How will we measure & test?
3. How will we scale?



Q/A

 ryan@emergetools.com

 @rbro112 / @emergetools