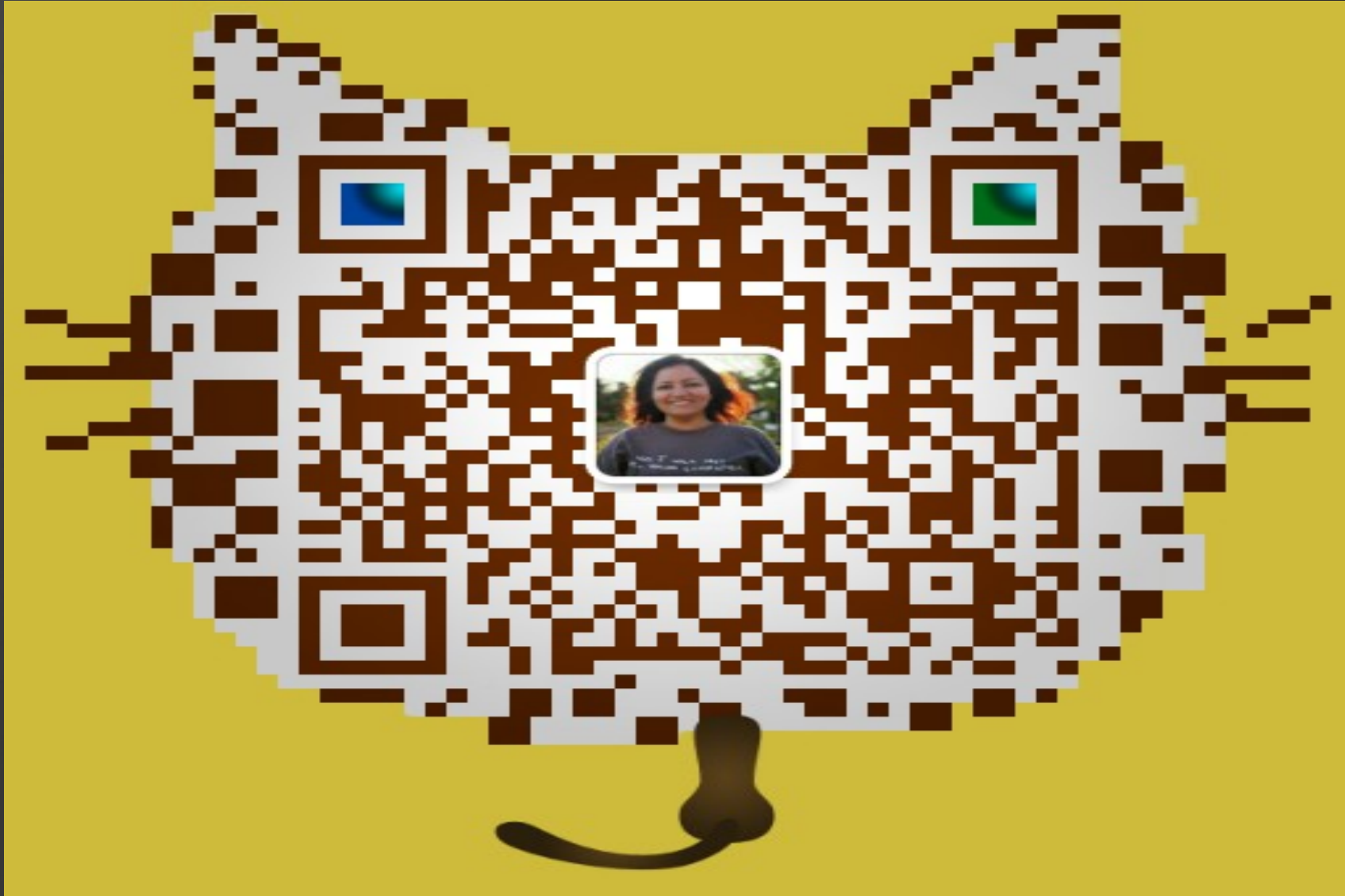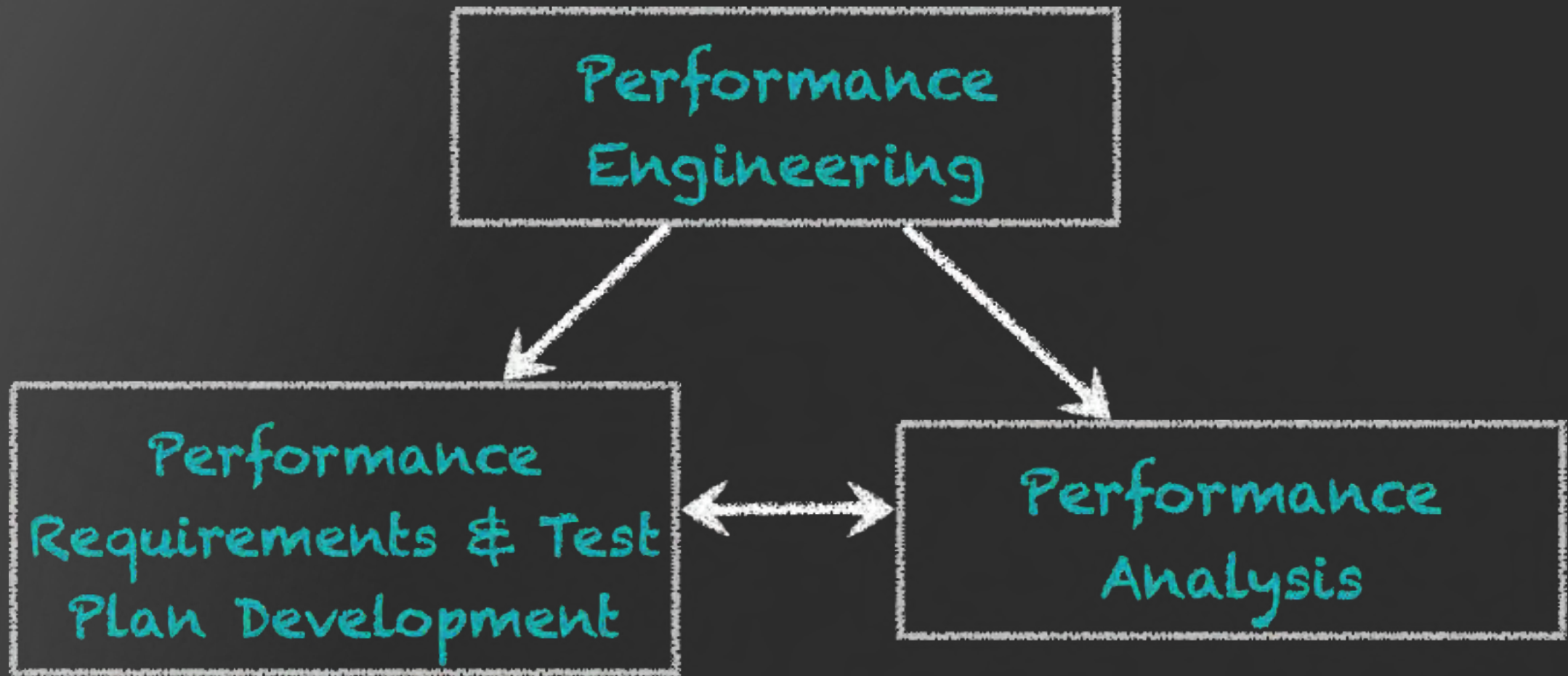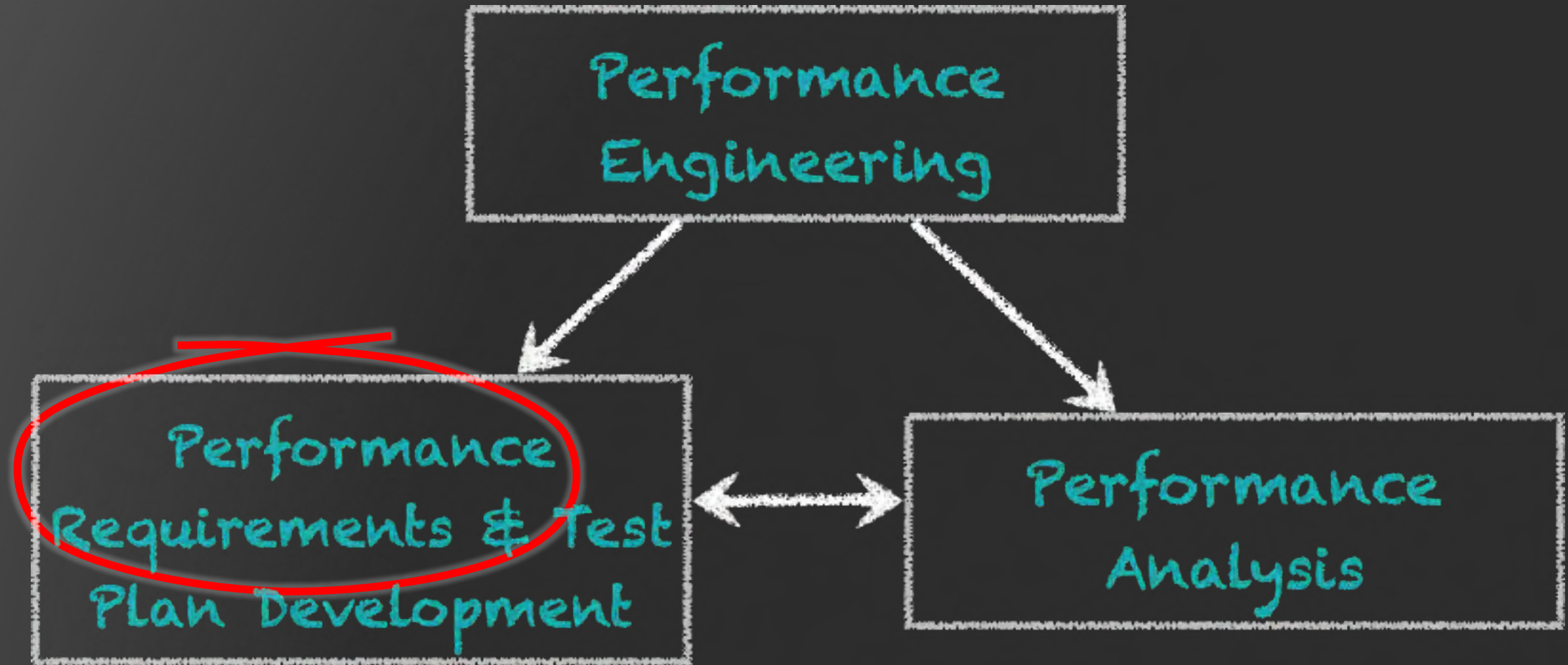# About Me

# Performance Engineering

# Performance Requirements

# Ask Questions

- What makes the users happy?

- What makes them unhappy?

- Any current issues that need to be tackled?

# Define and Understand

- Understand QoS (Quality of Service) for end users

- Define QoS success metrics in measurable terms

  - Those are the service level agreements (SLAs)!
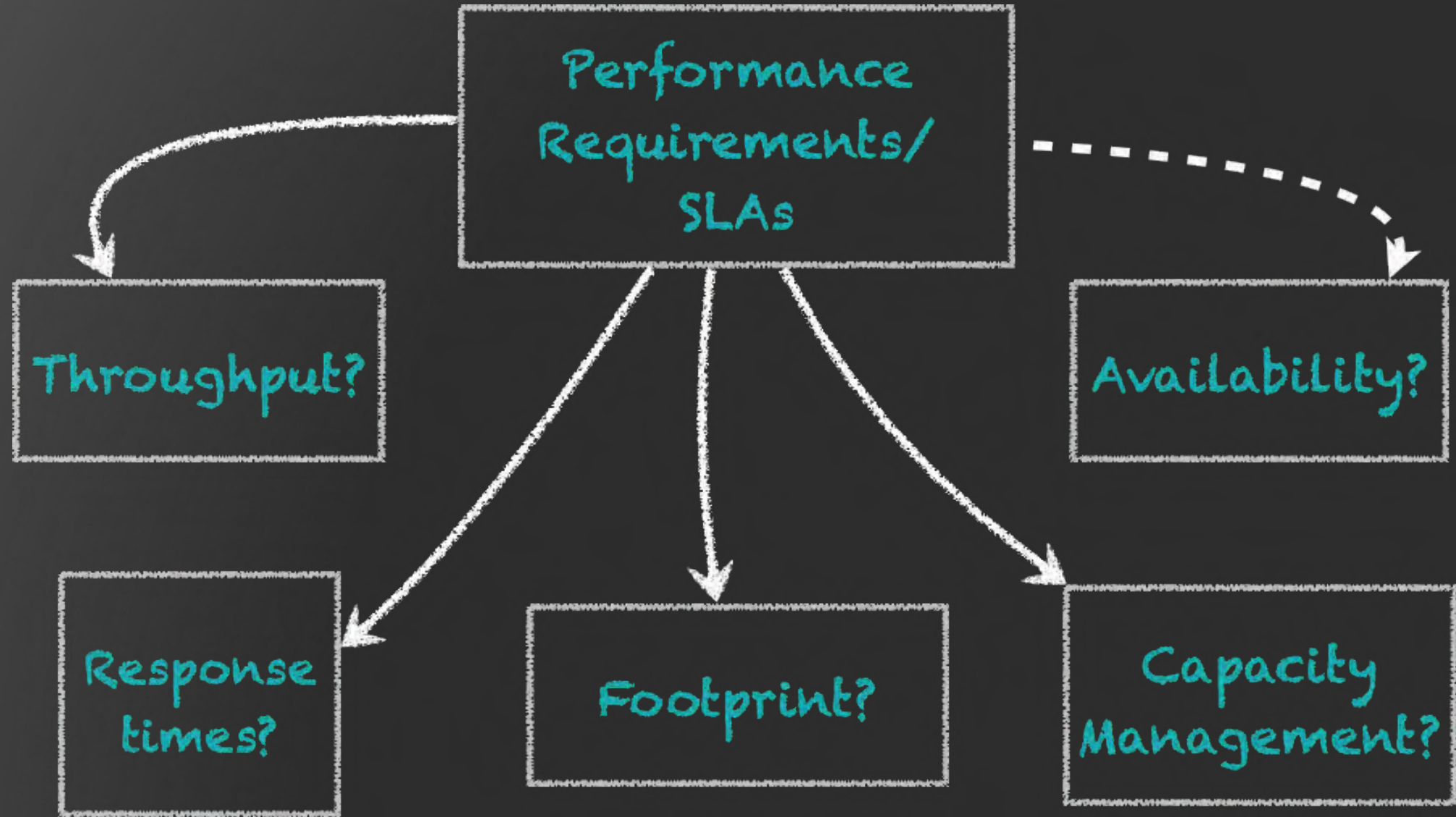
# Define Success!

- Document, understand and prioritize SLAs

  - Throughput

  - Response time

  - Capacity

  - Footprint
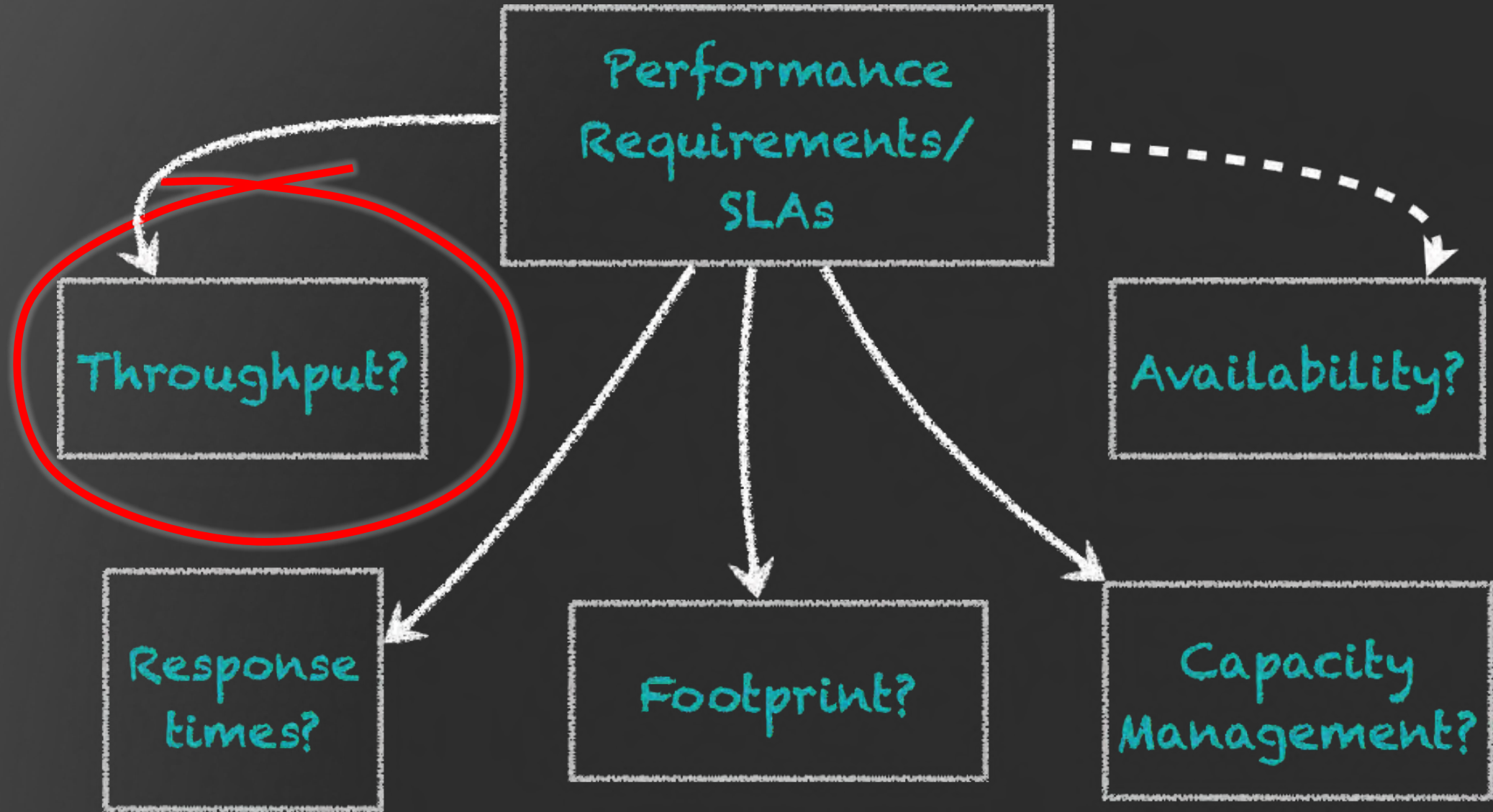
  - CPU utilization?

  - …

©2016 CodeKaram

# Work For It!

Monitor, measure and define performance in terms of throughput, latency, capacity, footprint, utilization ...

# Defining Success!

Performance Requirements/SLAs

Throughput?

Availability?

Response times?

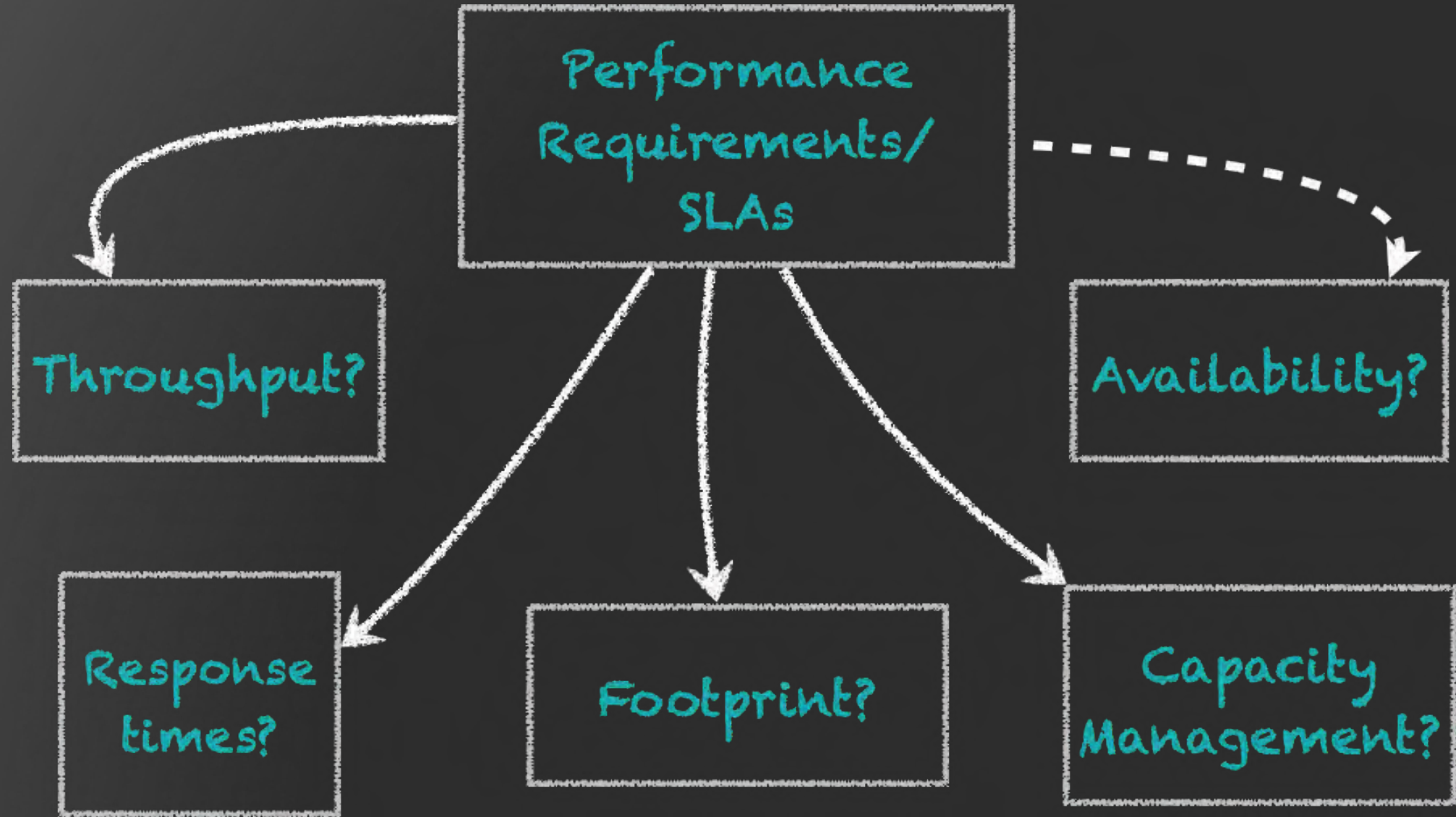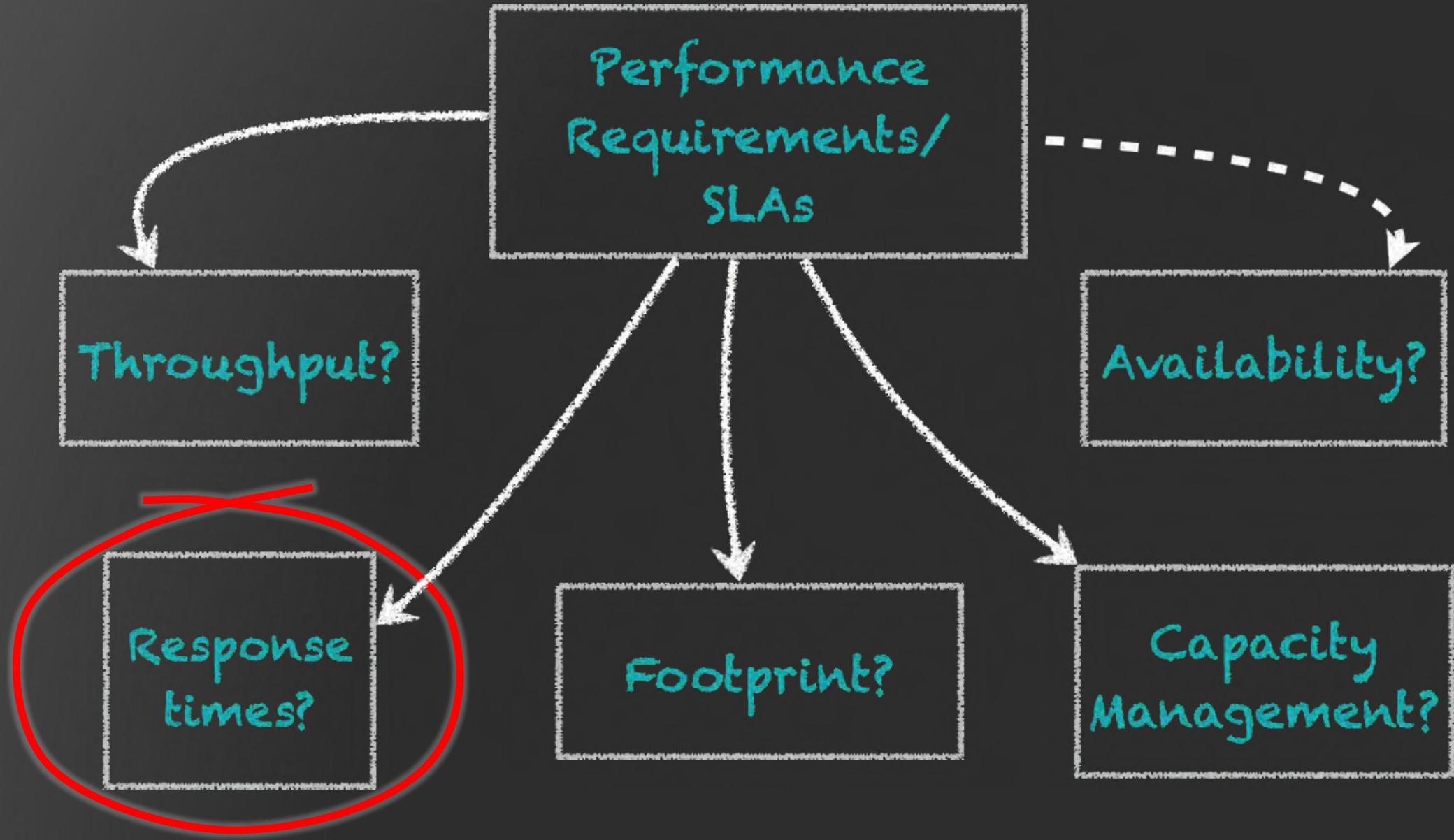Footprint?

Capacity Management?

# Throughput

- Expected throughput?

  - Can you fall below the expected throughput?

  - How long can you stay below the expected throughput?

  - What is the lowest that you can go?

# Throughput

- Throughput measurement

  - How is it measured?

    - Transactions/sec; messages/sec or all of them?

  - Where is it measured?

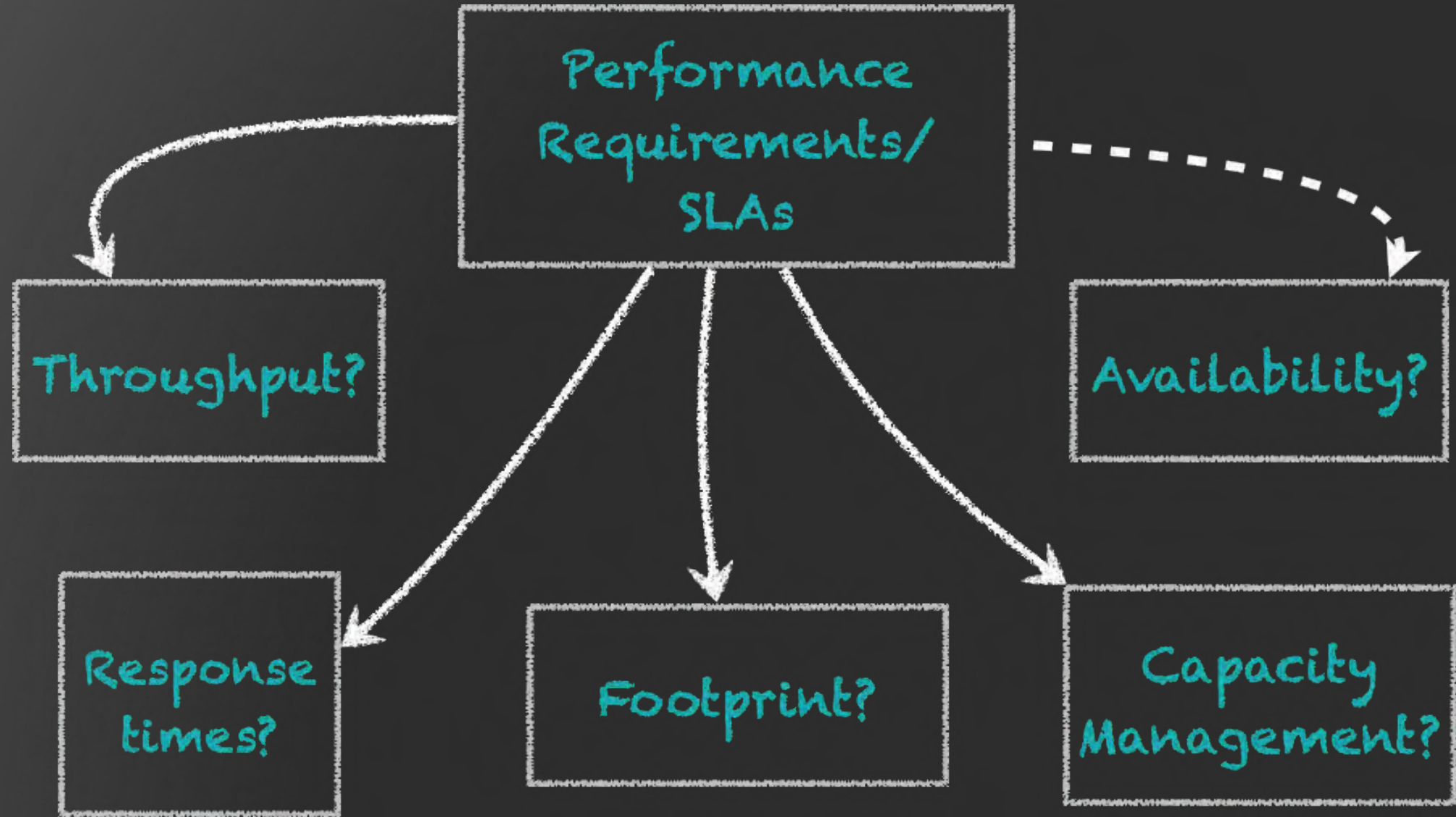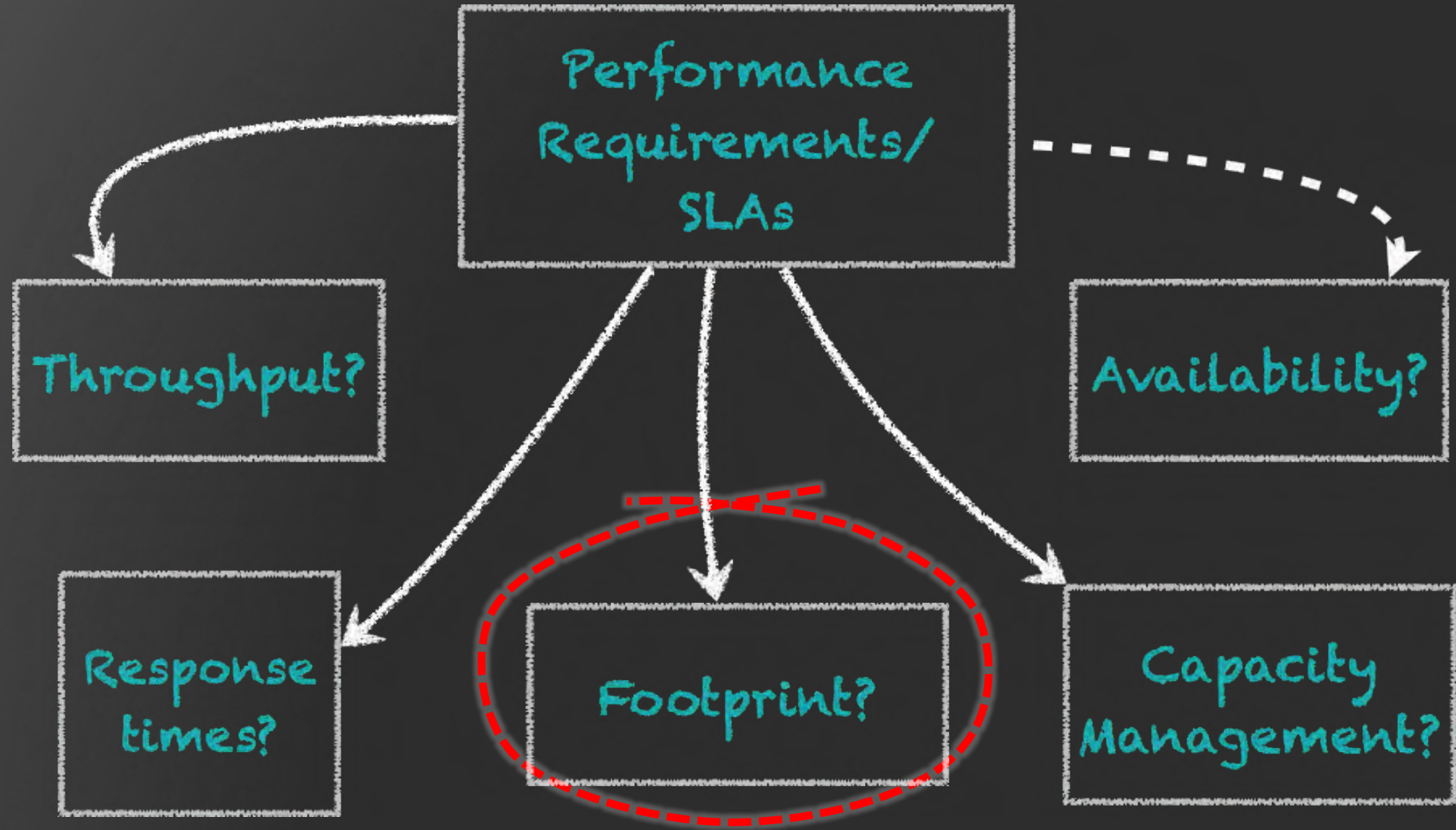    - Client-side; server-side; browser?

# Response Time

- Expected response time?

  - Can you go above the expected response time?

  - How long can you stay there?

  - How much can you exceed?

# Response Time

- Response time measurement

  - How is it measured?

    - 99$^{th}$ percentile; 5-9s; worst-case; or all?

  - Where is it measured?

    - Client-side; server-side? Complete loop?

# Capacity Management

- What is the expected load?

- What happens if one system gets loaded more than others? (load balancer issue)

- How is it measured?

# Capacity Management

- What's the max load that a system and all systems can tolerate?

- How long?

- What metrics are being captured?
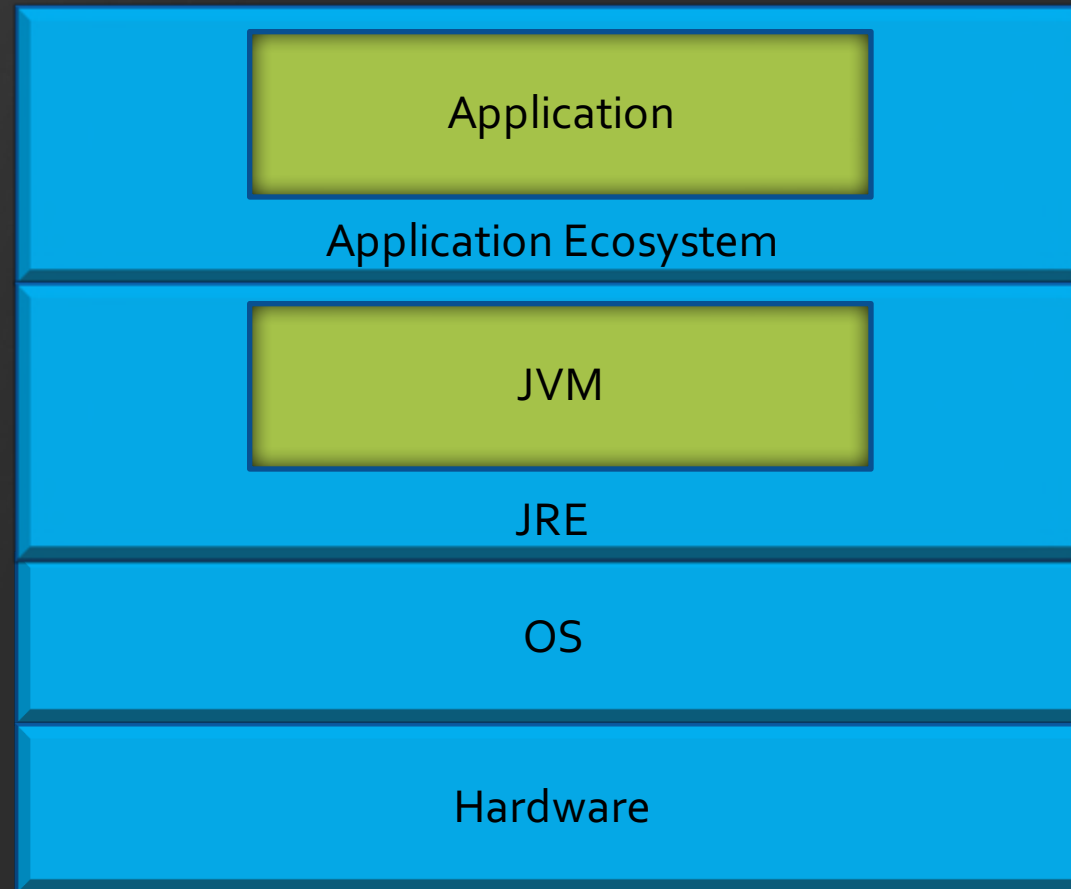
# Performance Analysis

# Performance Analysis

- Analyze what factors enable the end-user experience to meet or exceed the promised QoS
- Track your SLAs!

# Java Application Stack

Application

Application Ecosystem

JVM

JRE

OS

Hardware

# Application Performance Analysis



- Application services
- Application server
- Database
- Any other services in the ecosystem?

# JVM + Runtime Performance Analysis



- Classloading stats
- JIT Compilation stats
- Garbage Collection stats
- Threads stats

# OS Performance Analysis

OS

- System/ Kernel stats
- Lock stats
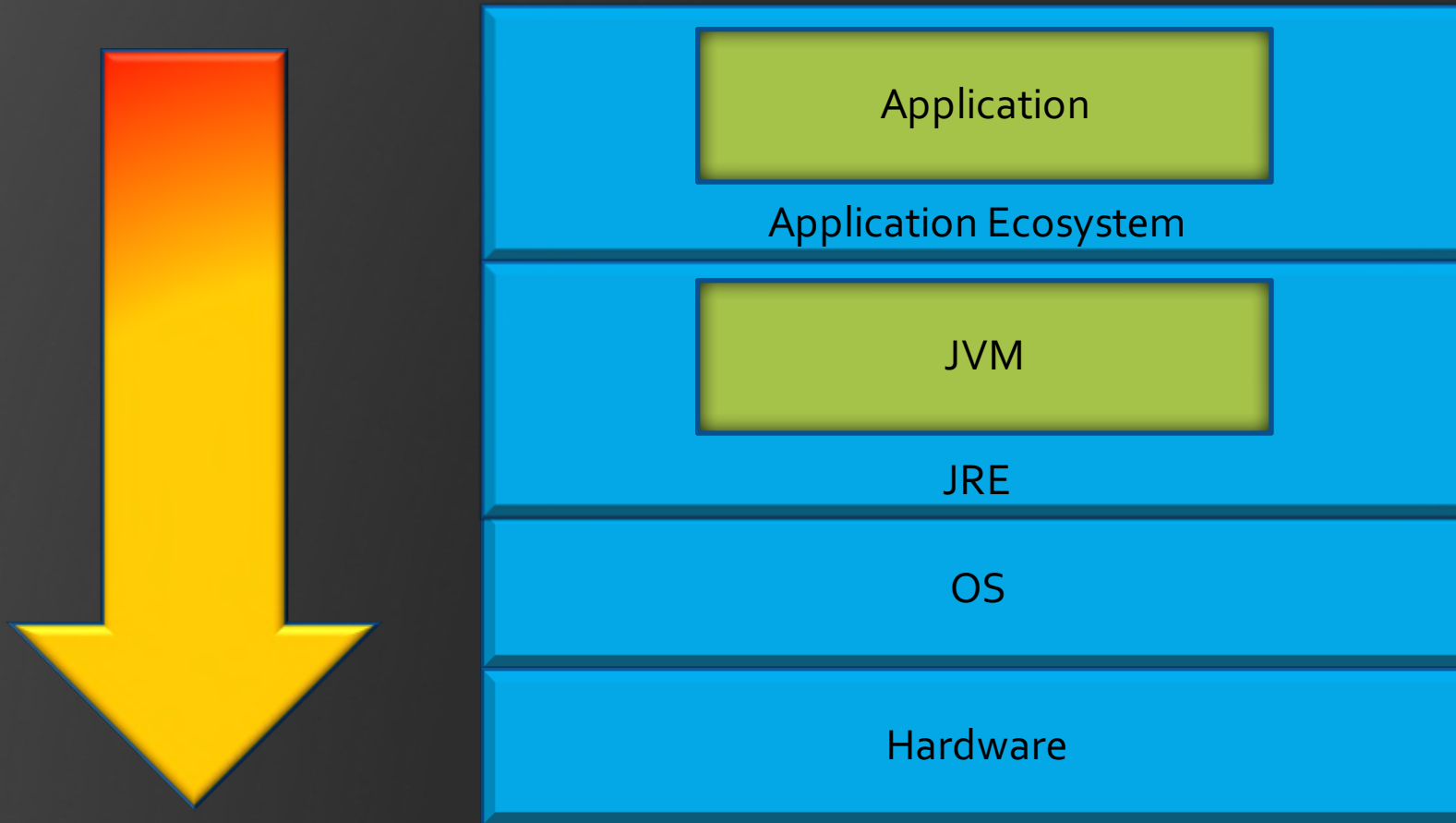- Threads stats

# Hardware Performance Analysis

Hardware

- Memory bandwidth/ traffic/ consumption
- CPU/ core utilization
- CPU cache efficiency/ utilization/ levels
- Architectural specific?
- IO Stats

# What are You Trying to Achieve?

Improve application?

# Top-Down Approach

Application

Application Ecosystem

JVM

JRE

OS

Hardware

# Top Down Approach - Process



Monitor • Step 1

• Step 2 Profile

Analyze • Step 3

• Step 4 Tune + Apply

# What are You Trying to Achieve?

Improve the platform?

# Bottom-Up Approach

Application

Application Ecosystem

JVM

JRE

OS

Hardware

# Bottom Up Approach - Process



**Monitor** • **Step 1**

• **Step 2** **Profile**

**Analyze** • **Step 3**

• **Step 4** **Tune + Apply**

# Top-Down Approach

# Top-Down Approach

I HAVE the power!!

... to modify the code

the code

# Top-Down Approach



Application

Application Ecosystem

JVM

JRE

OS

Hardware

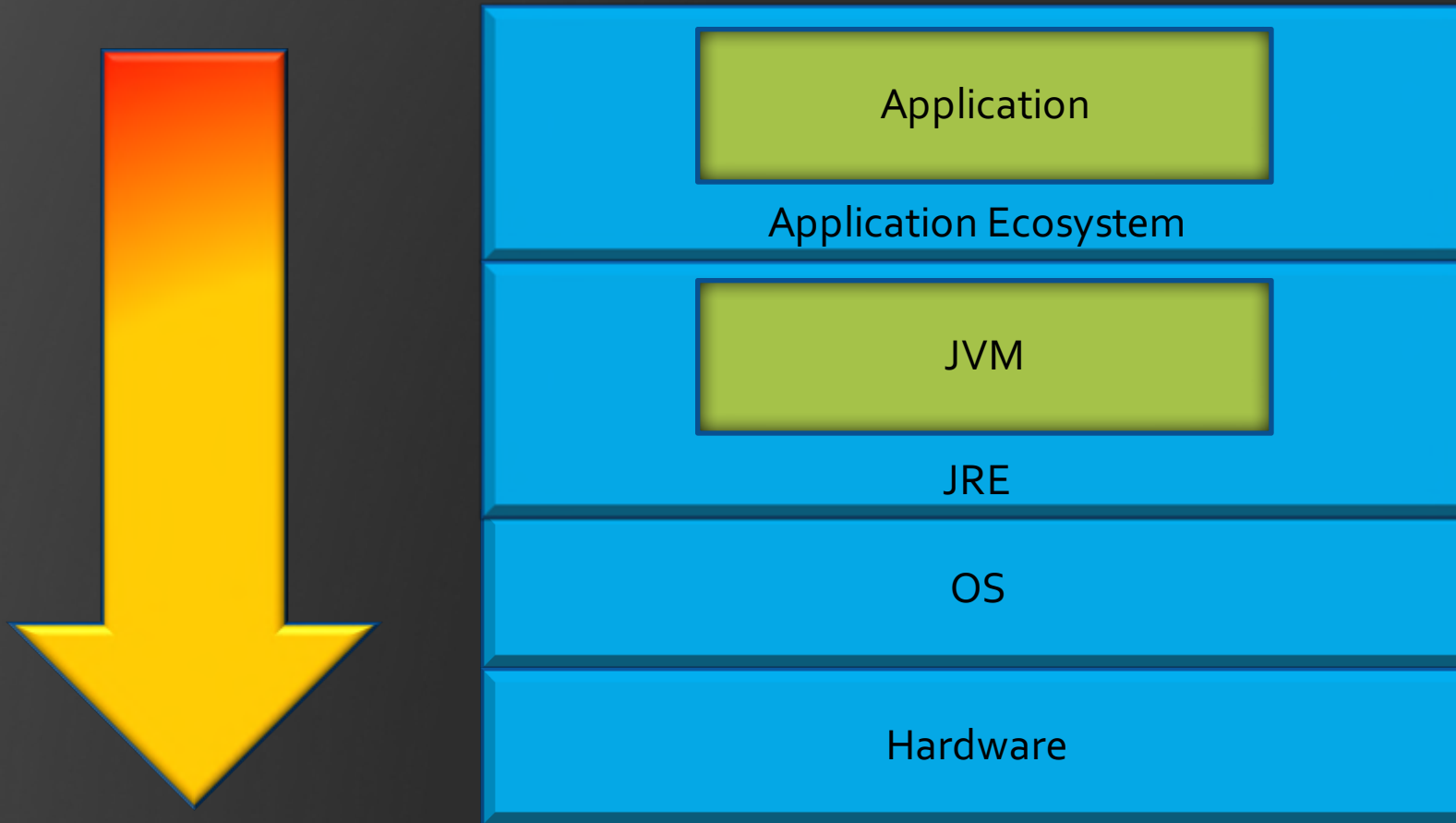# Top-Down Approach

Application

Application Ecosystem

JVM

JRE

**Platform**

OS

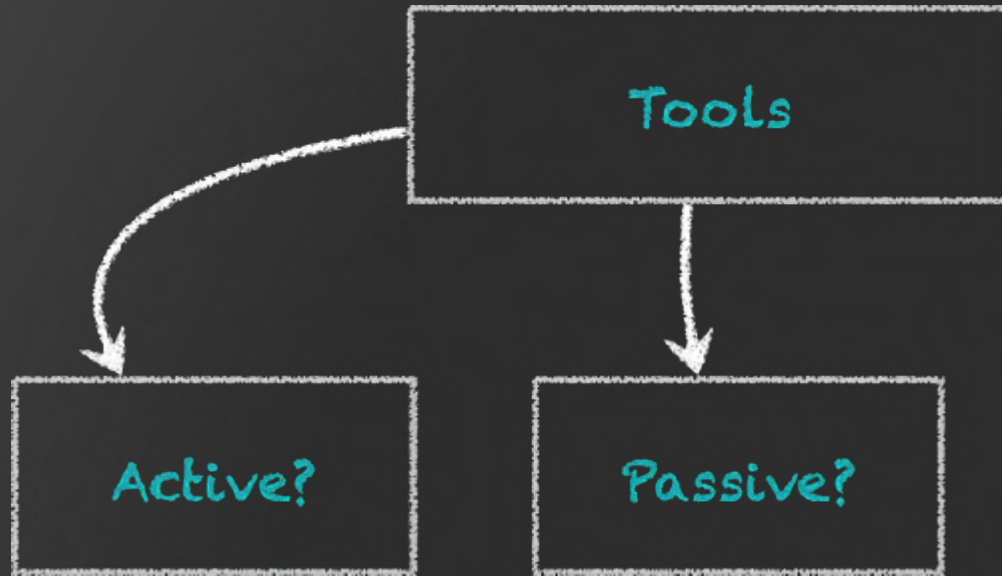Hardware

# Step 1: Monitor

Monitor and plot SUT (System Under Test) statistics

# Step 1: Monitor

# Step 1: Monitor

```
                    ┌─────────────┐
                    │    Tools    │
                    └─────────────┘
                  ↙                ↓
        ┌─────────────┐      ┌─────────────┐
        │   Active?   │      │  Passive?   │
        └─────────────┘      └─────────────┘
              ↓
        ┌─────────────┐
        │   Create    │
        │   alerts    │
        └─────────────┘
```

# Step 1: Monitor

# Step 1: Monitor

# Step 1: Monitor



Tools

Active? → Create alerts

Passive? → Network taps

Offline? → Capture logs

# Step 1: Monitor

- Tools?

  - VisualVM, Java Flight Recorder

  - PrintCompilation, PrintGCDetails (+PrintGCDateStamps), jmap – clstats, jcmd GC.class_stats

# Step 1: Monitor

- Tools?
  - Linux – mpstat, sysstat – iostat,pidstat…, prstat, vmstat, dash, CPU-Z, cacti …
  - Windows – Performance Monitor, Task Manager, Resource Monitor, CPU-Z, cacti …

# Step 2+3: Profile + Analyze

- You have all the data that you need!
- Identify areas of improvement
- Profile those potential hotspots
- Analyze those hotspots

# Step 2+3: Profile + Analyze

- Tools? (Free/Open source/GPL/BSD)

  - Oracle Solaris Studio Performance Analyzer, perf tools, PAPI, Code XL, Dtrace, Oprofile, gprof, LTT (linux trace toolkit)

  - Java Application – VisualVM, Netbeans profiler, jconsole …

# Step 4: Tune

- Tune the JVM/GC – select the right heap, the right GC algorithm
  - Age objects appropriately
  - Promote only long-lived objects
  - GC worker threads per VM (for stop-the-world GC events)
  - GC concurrent worker threads per VM

# Step 4: Tune

- Tune the JVM/GC – select the right heap, the right GC algorithm
  - See if compressed oops can be helpful
  - Larger heaps may need AlwaysPretouch to be enabled and also UseLargePages of appropriate size

# Step 4: Tune

- Tune your code to meet or exceed your SLAs
  - Appropriate ramp-ups and ramp-downs
  - Object age and retention strategies
    - Understand what forms your LDS (live data set)
  - Confirm you are measuring the right thing! ☺
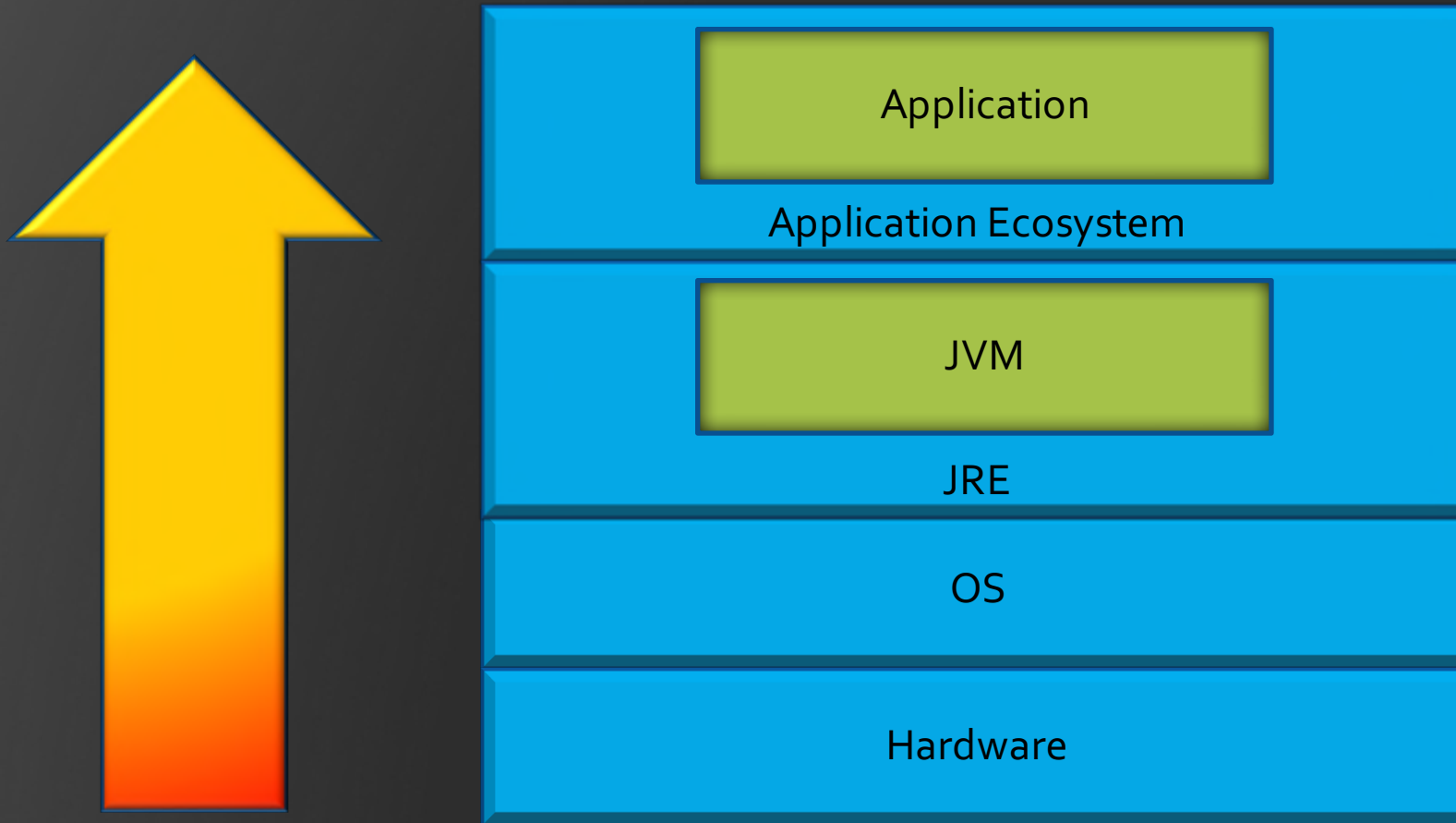
# Bottom-Up Approach
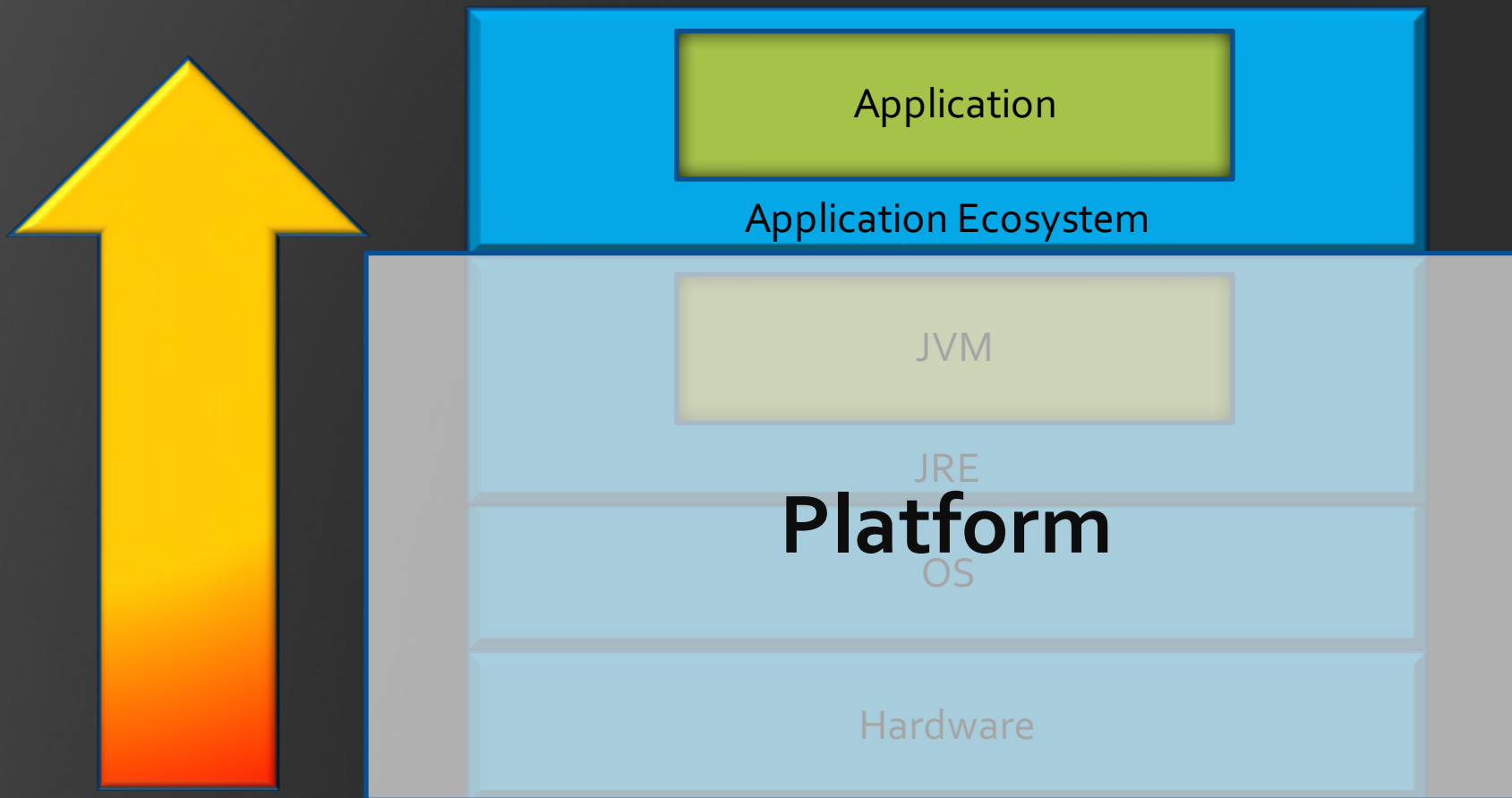
# Bottom Up Approach

I NEED the power!!

... to stress the platform

# Top-Down Approach

# Top-Down Approach

Application

Application Ecosystem

JVM

JRE

**Platform**

OS

Hardware

# Where to Start?

- Know what you are stressing

- Get/ write the appropriate workload/ application

- Get/write the appropriate tools

©2016 CodeKaram

# Know What You Are Stressing!

CPU –
Gather performance counter information for your CPU stats, core stats, cache hits, misses and levels, branch predictions, pipeline information, order-of-execution, load-store unit load and queues, etc

©2016 CodeKaram

# Know What You Are Stressing!

Memory –
Gather performance counter information for memory utilization, memory bandwidth, read-write stats, max read bandwidth, max write bandwidth, max cross traffic bandwidth, architectural related considerations, max capacity, etc

# Know What You Are Stressing!

JVM / GC–
Gather information related to the change – e.g. new GC!
Gather information on different GC phases, parallel work queues and work performance, concurrent work and pressure, internal queues and buffers, any GC work that's staged?, etc

# Know What You Are Stressing!

JVM / GC–

Gather information related to the change – e.g. new GC!

Gather information on different GC phases, parallel work queues and work performance, concurrent work and pressure, internal queues and buffers, any GC work that's staged?, etc

# Where to Next?

- Know what you are stressing ✓

- Get/ write the appropriate workload/ application ✓

- Get/write the appropriate tools ✓

Let's have some fun!

# Further Reading

**Java性能优化圣经！Java之父重磅推荐！**

# Java 性能优化权威指南

Java Performance

[美] Charlie Hunt
Binu John ■ 著

柳飞 陆明刚 ■ 译

www.codekaram.com
www.slideshare.net/monicabeckwith