# Research Methods

CSCI 8901:
Visualizing Your Research

Prof. Tim Wood
GWU
2021

# Visuals Matter

Slides, Videos, Posters, Demos

Papers
- Diagrams
- Graphs
- Even fonts and formatting!

You want your work to look:
- Professional
- Attractive
- Memorable
- Informative

# Let's go back in time…

Tim Wood's Thesis Defense

Sometime in April, 2011
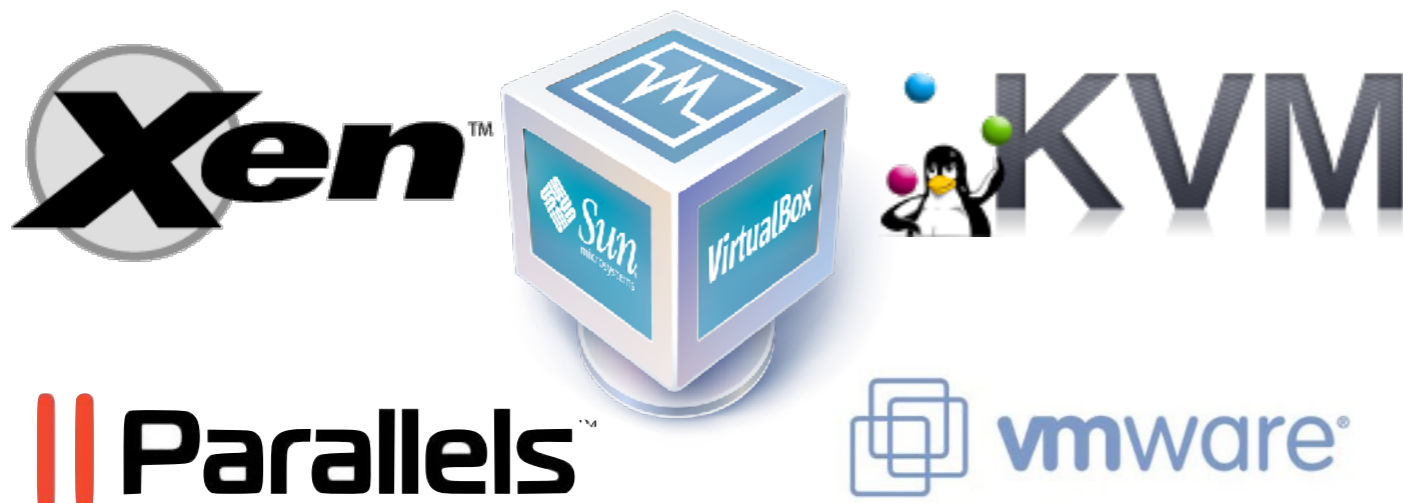Somewhere in farm country, Massachusetts

# Data Centers

- **Infrastructure as a Service** clouds rent server and storage resources on demand

- Data Centers are BIG server farms
  - Clusters of 10,000s of servers
  - Growing to 100s of thousands

- Host many application types
  - Web servers, databases
  - Custom business apps
  - Search clusters, data mining

Challenges: large scale and dynamic workload fluctuations

# Server Virtualization

- Data centers use virtualization to share physical resources and simplify automation

- Allows a server to be "sliced" into Virtual Machines

- VM has own OS/applications

- Rapidly adjust resource allocations

- VM migration within a LAN

# Within a Data Center

- How to transition applications to VMs and account for virtualization overheads?

    - **MOVE**: Modeling Overheads of Virtual Environments

- Where should VMs be placed to allow for the greatest level of server consolidation?

    - **Memory Buddies**: Memory sharing guided placement

- How to dynamically allocate VM resources to prevent server overload?

    - **Sandpiper**: Automated VM migration and resizing
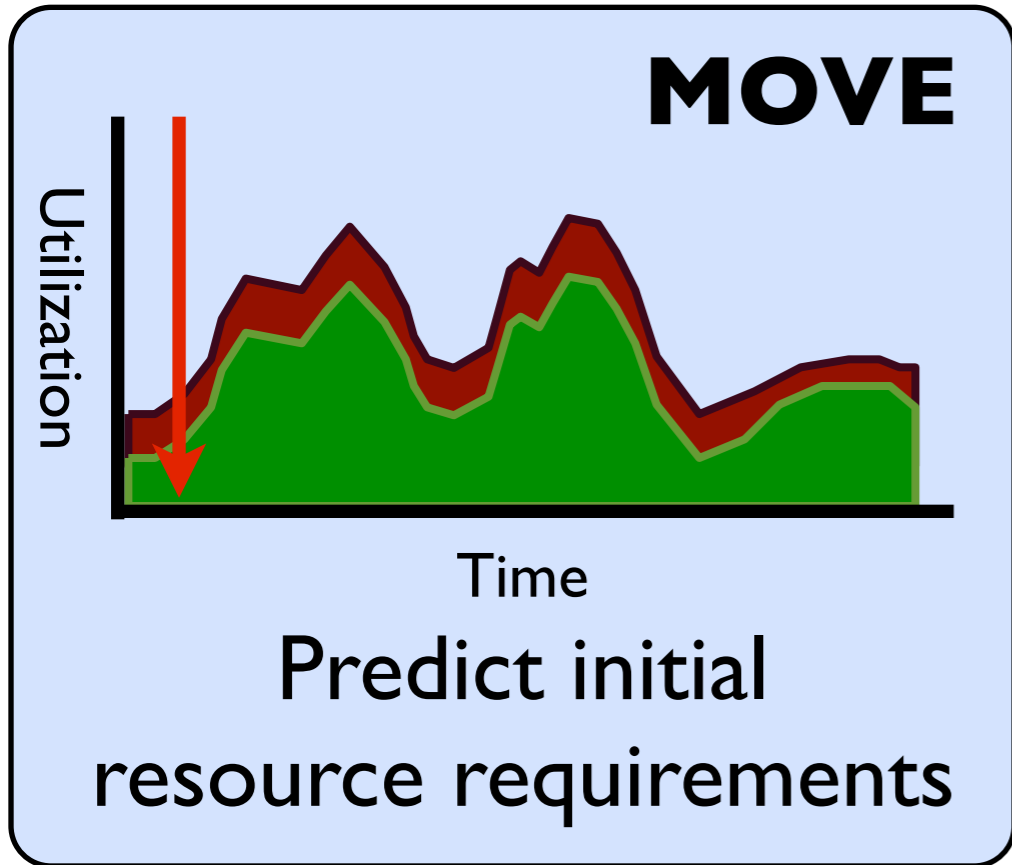
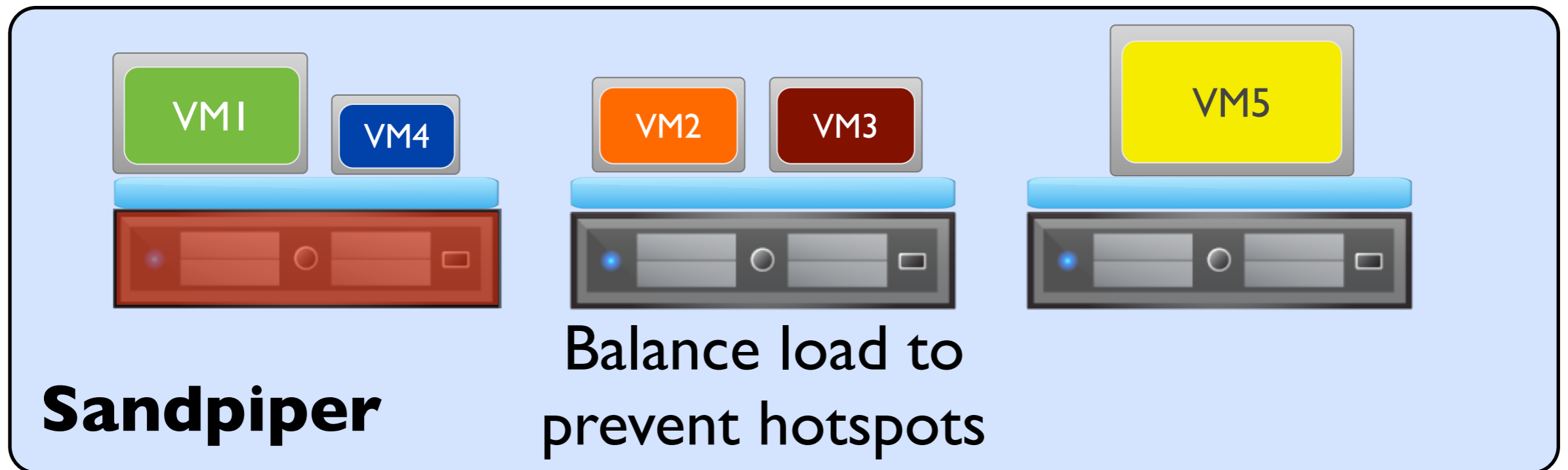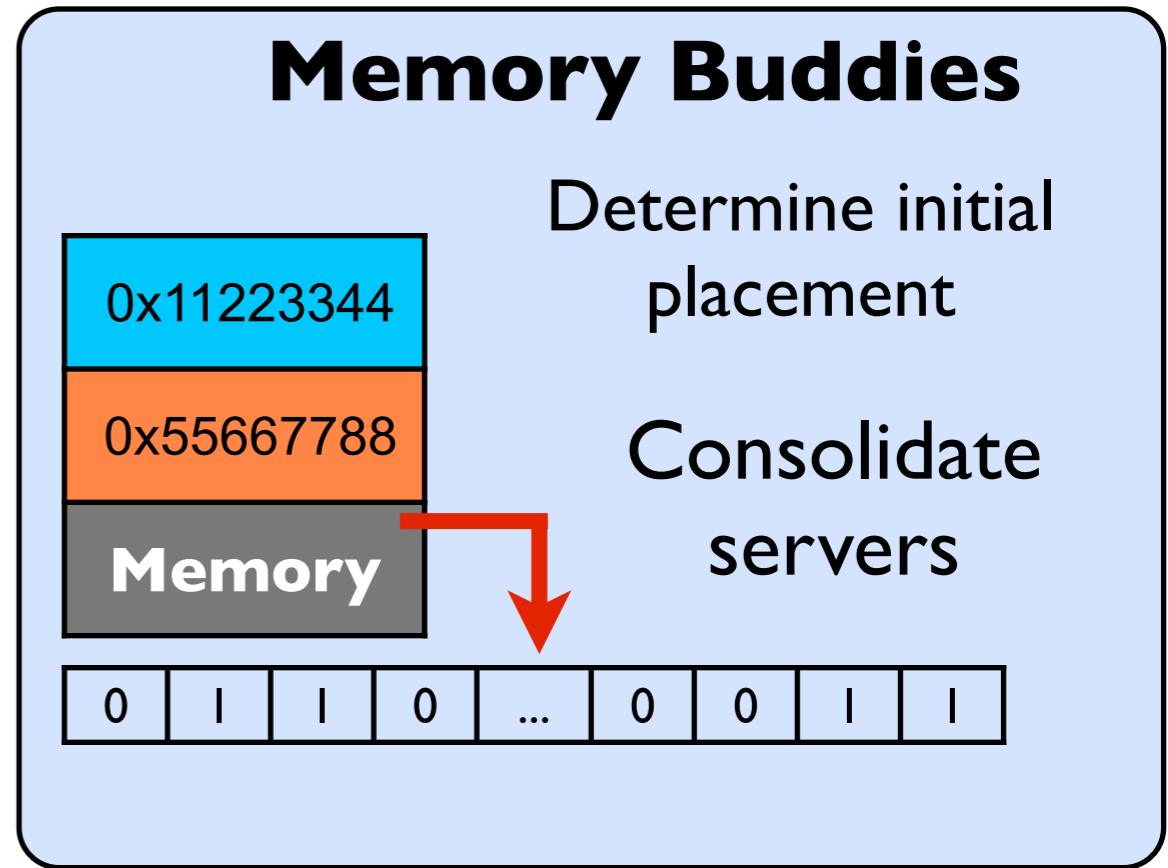**Deployment**  **Resource Management**  **Reliability**

| MOVE | Memory Buddies | Sandpiper | CloudNet | Pipe Cloud |
|------|----------------|-----------|----------|------------|

# Within a Data Center

## MOVE



Utilization

Time

Predict initial resource requirements

VM1
VM2
VM3
VM4
VM5

## Memory Buddies

Determine initial placement

Consolidate servers

0x11223344
0x55667788
**Memory**

| 0 | 1 | 1 | 0 | ... | 0 | 0 | 1 | 1 |

## Sandpiper

VM1  VM4

VM2  VM3

VM5

Balance load to prevent hotspots

# What is good in these slides?

(or bad!)

(your ideas here)

# What is good in these slides?

Good balance of text and visuals

Nice high level overview of thesis

Good connections between components of thesis


Animation is a bit excessive

# Slide Tips

# US Wireless Market – Q2 2010 Update

## What to expect in the coming months?

31% of the US subscription base is now smartphones.

The pace of product introduction is accelerating with each quarter. Devices of all shapes and sizes are coming into the market literally every week. Players are having to re-evaluate their businesses and long-term strategies. Several new impressive devices got introduced during the course of 1H of 2010 including the iPad and EVO.

There are several players whose future is at stake (to put it mildly). The competition has grown fierce and companies are finding it hard to take ideas from R&D to products in market in a short amount of time.

Microsoft announced its comeback with the W7 commercial launch imminent. The change in UI was refreshing and the expectations are quite high. W7 v2 is likely around the corner to update on the flaws of v1. HP acquired Palm in an attempt to become relevant again in the mobile device space. It has been an action packed 1H 2010 and we can expect more of the same for the remainder of the year.

2010 has also been active on the regulatory front as the national broadband plan was unveiled in March and the subsequent debate over the course of nations broadband future kept the spectrum, net-neutrality, and exclusivity issues at the forefront.

To start planning for 4G, 5G, and beyond, US should think about rolling a 50 year broadband plan. While more spectrum is always helpful, will we have all the spectrum we need in 2050? or do we need to invent new technologies and business models that use spectrum more wisely? This topic will keep the industry occupied for some time to come. (We will be going in-depth into this subject at our Sept event with some very senior and experienced executives)

2010 is also the year of network rollouts. T-Mobile has been rolling out HSPA+ at an impressive rate, Clearwire announced its intention to move to LTE, Verizon is betting big on LTE and looking for competitive marketing advantage over the course of the next 12 months. AT&T has been adding backhaul, upgrading to HSPA+ and planning for LTE all at once. Even the smaller carriers like MetroPCS are looking for competitive advantage with quicker LTE launch and beat others by carrying the first LTE smartphone. (We will be releasing the next edition of our "State of the "Mobile" Broadband Nation" paper later this year)

As we had mentioned last year, the mobile data traffic kept on growing disproportional to the revenues. A series of solutions have come into the market from players big and small. We released the second edition of our in-depth research paper on data growth - "Managing Growth and Profits in the Yottabyte Era" last quarter.

We will be keeping a very close eye on the micro- and macro-trends and reporting on the market on a regular basis in various private and public settings.

# Limit your text (84 point)

Use large fonts (41 point)
- Not smaller than this (32 point)

Use bullets, not paragraphs
- **Emphasize** your key points

Don't try to be exhaustive
- Unless the slides will be referred to later without your speech

Don't try to cram in too much content!

## US Wireless Market – Q2 2010 Update

**What to expect in the coming months?**

31% of the US subscription base is now smartphones.

The pace of product introduction is accelerating with each quarter. Devices of all shapes and sizes are coming into the market literally every week. Players are having to re-evaluate their businesses and long-term strategies. Several new impressive devices got introduced during the course of 1H of 2010 including the iPad and EVO.

There are several players whose future is at stake (to put it mildly). The competition has grown fierce and companies are finding it hard to take ideas from R&D to products in market in a short amount of time.

Microsoft announced its comeback with the W7 commercial launch imminent. The change in UI was refreshing and the expectations are quite high. W7 v2 is likely around the corner to update on the flaws of v1. HP acquired Palm in an attempt to become relevant again in the mobile device space. It has been an action packed 1H 2010 and we can expect more of the same for the remainder of the year.

2010 has also been active on the regulatory front as the national broadband plan was unveiled in March and the subsequent debate over the course of nations broadband future kept the spectrum, net-neutrality, and exclusivity issues at the forefront.

To start planning for 4G, 5G, and beyond, US should think about rolling a 50 year broadband plan. While more spectrum is always helpful, will we have all the spectrum we need in 2050? or do we need to invent new technologies and business models that use spectrum more wisely? This topic will keep the industry occupied for some time to come. (We will be going in-depth into this subject at our Sept event with some very senior and experienced executives)

2010 is also the year of network rollouts. T-Mobile has been rolling out HSPA+ at an impressive rate, Clearwire announced its intention to move to LTE, Verizon is betting big on LTE and looking for competitive marketing advantage over the course of the next 12 months. AT&T has been adding backhaul, upgrading to HSPA+ and planning for LTE all at once. Even the smaller carriers like MetroPCS are looking for competitive advantage with quicker LTE launch and beat others by carrying the first LTE smartphone. (We will be releasing the next edition of our "State of the "Mobile" Broadband Nation" paper later this year)

As we had mentioned last year, the mobile data traffic kept on growing disproportional to the revenues. A series of solutions have come into the market from players big and small. We released the second edition of our in-depth research paper on data growth - "Managing Growth and Profits in the Yottabyte Era" last quarter.

We will be keeping a very close eye on the micro- and macro-trends and reporting on the market on a regular basis in various private and public settings.

# Limit your text (ugly version)

Use large fonts

- Not smaller than this (32 point)

Use bullets, not paragraphs

- Emphasize your key points

Don't try to be exhaustive

- Unless the slides will be referred to later without your speech

Don't try to cram in too much content!



**US Wireless Market – Q2 2010 Update**

What to expect in the coming months?

31% of the US subscription base is now smartphones.

The pace of product introduction is accelerating with each quarter. Devices of all shapes and sizes are coming into the market literally every week. Players are having to re-evaluate their businesses and long-term strategies. Several new impressive devices got introduced during the course of 1H of 2010 including the iPad and EVO.

There are several players whose future is at stake (to put it mildly). The competition has grown fierce and companies are finding it hard to take ideas from R&D to products in market in a short amount of time.

Microsoft announced its comeback with the W7 commercial launch imminent. The change in UI was refreshing and the expectations are quite high. W7 v2 is likely around the corner to update on the flaws of v1. HP acquired Palm in an attempt to become relevant again in the mobile device space. It has been an action packed 1H 2010 and we can expect more of the same for the remainder of the year.

2010 has also been active on the regulatory front as the national broadband plan was unveiled in March and the subsequent debate over the course of nations broadband future kept the spectrum, net-neutrality, and exclusivity issues at the forefront.

To start planning for 4G, 5G, and beyond, US should think about rolling a 50 year broadband plan. While more spectrum is always helpful, will we have all the spectrum we need in 2050? or do we need to invent new technologies and business models that use spectrum more wisely? This topic will keep the industry occupied for some time to come. (We will be going in-depth into this subject at our Sept event with some very senior and experienced executives)

2010 is also the year of network rollouts. T-Mobile has been rolling out HSPA+ at an impressive rate, Clearwire announced its intention to move to LTE, Verizon is betting big on LTE and looking for competitive marketing advantage over the course of the next 12 months. AT&T has been adding backhaul, upgrading to HSPA+ and planning for LTE all at once. Even the smaller carriers like MetroPCS are looking for competitive advantage with quicker LTE launch and beat others by carrying the first LTE smartphone. (We will be releasing the next edition of our "State of the "Mobile" Broadband Nation" paper later this year)

As we had mentioned last year, the mobile data traffic kept on growing disproportional to the revenues. A series of solutions have come into the market from players big and small. We released the second edition of our in-depth research paper on data growth - "Managing Growth and Profits in the Yottabyte Era" last quarter.

We will be keeping a very close eye on the micro- and macro-trends and reporting on the market on a regular basis in various private and public settings.

# Mix Text and Images

# USE A HIERARCHY

- This is text

  - This is also text

    - This is even more text

- Why are they all the same size?

# Use a hierarchy

This is text
- This is also text
    - This is even more text
- Note that they are not the same size!

White space is important, but don't go overboard
- And sub bullets with a smaller font size help viewers focus on key points

Make your own template and keep improving it!

## USE A HIERARCHY

- This is text

    - This is also text

        - This is even more text

    - Why are they all the same size?

# GW PPT Template

This is where I put my content
            Here is more content
Wow, this is just awful.
Why is the bar so big at the bottom?


I have so little useful space and it is poorly laid out.

# My Template

Text that is reasonably large

- Sub bullets that are smaller
    - Sub-sub bullets that are even smaller, although I rarely use them
    - (Mainly so I can add spacing more flexibly)

Large "before paragraph" spacing so bullets aren't too tight and smaller line spacing so you can fit denser text when needed (try to avoid multi-line)

A useful footer with your name and affiliation

- Always include the slide number in corner!

Minimal background images

Optional: school / lab logos

**Pop up boxes to emphasize key points!**

# Should we animate?

Text that is reasonably large
  - Sub bullets that are smaller
    - Sub-sub bullets that are even smaller, although I rarely use them
    - (Mainly so I can add spacing more flexibly)

Large "before paragraph" spacing so bullets aren't too tight and smaller line spacing so you can fit denser text when needed (try to avoid multi-line)

A useful footer with your name and affiliation
  - Always include the slide number in corner!
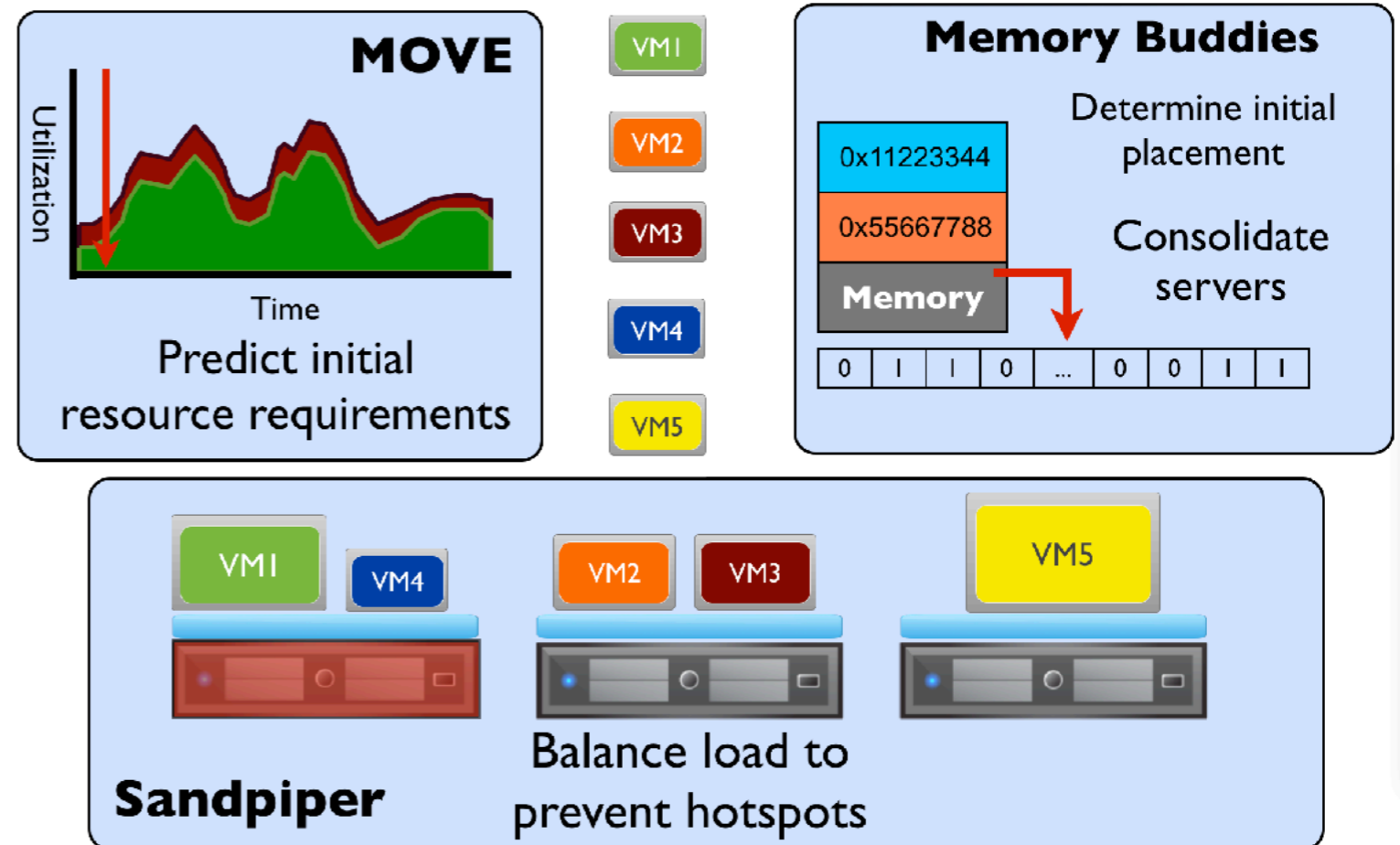
Minimal background images

Optional: school / lab logos

**It depends! Usually only if ~3 bullets on slide**

# Animations

Useful, but "expensive" to create

Can be distracting if overused



**Suggestion**: only use animation for emphasizing most important points

- And you only have at most 3 of those, right?

# Color Inspiration

Fig 2

Fig 3

Fig 4

Fig 7

From Boxes and Arrows

# Color

## Related colors


©Jill Morton - Color Matters

## Complementary colors


©Jill Morton - Color Matters

Limit the number of colors

Max per display: 4

Max across entire app: 7

# Font/Background Color

White background with a black font is easier to read

Black background with white font can look childish

Other colors may not have enough contrast or could look strange depending on the projector

# Font/Background Color

White background with a black font is easier to read

Black background with white font can look childish

Other colors may not have enough contrast or could look strange depending on the projector

# Font/Background Color

White background with a black font is easier to read

Black background with white font can look childish

Other colors may not have enough contrast or could look strange depending on the projector

# Fonts

Know the difference between:

Serif fonts: easier to read in print

- Times New Roman

Sans-Serif fonts: more modern on screen

- **Arial**, Helvetica

Monospaced fonts: only for code

- `Courier`

**Never use Comic Sans!**

# Ted Talk style?

Should we mimic TED talk slide style?

PREPARING TO FAIL
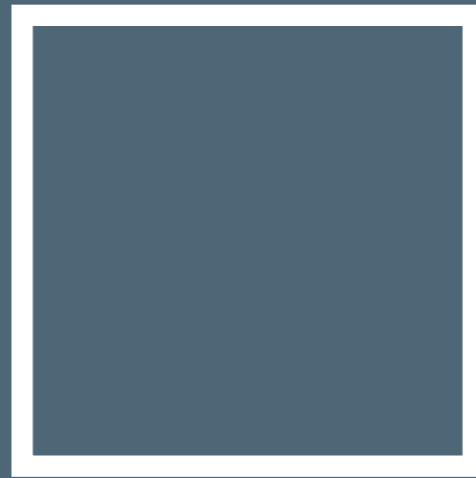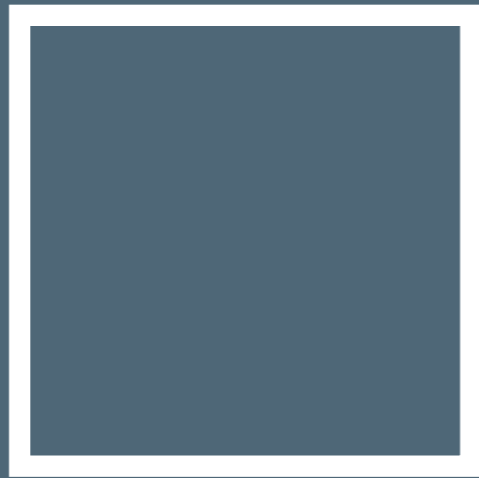
Photo: Blair Harkness

# PEOPLE

---

# CODE

---

# INFRASTRUCTURE

# INFRASTRUCTURE

"Success is stumbling from failure to failure
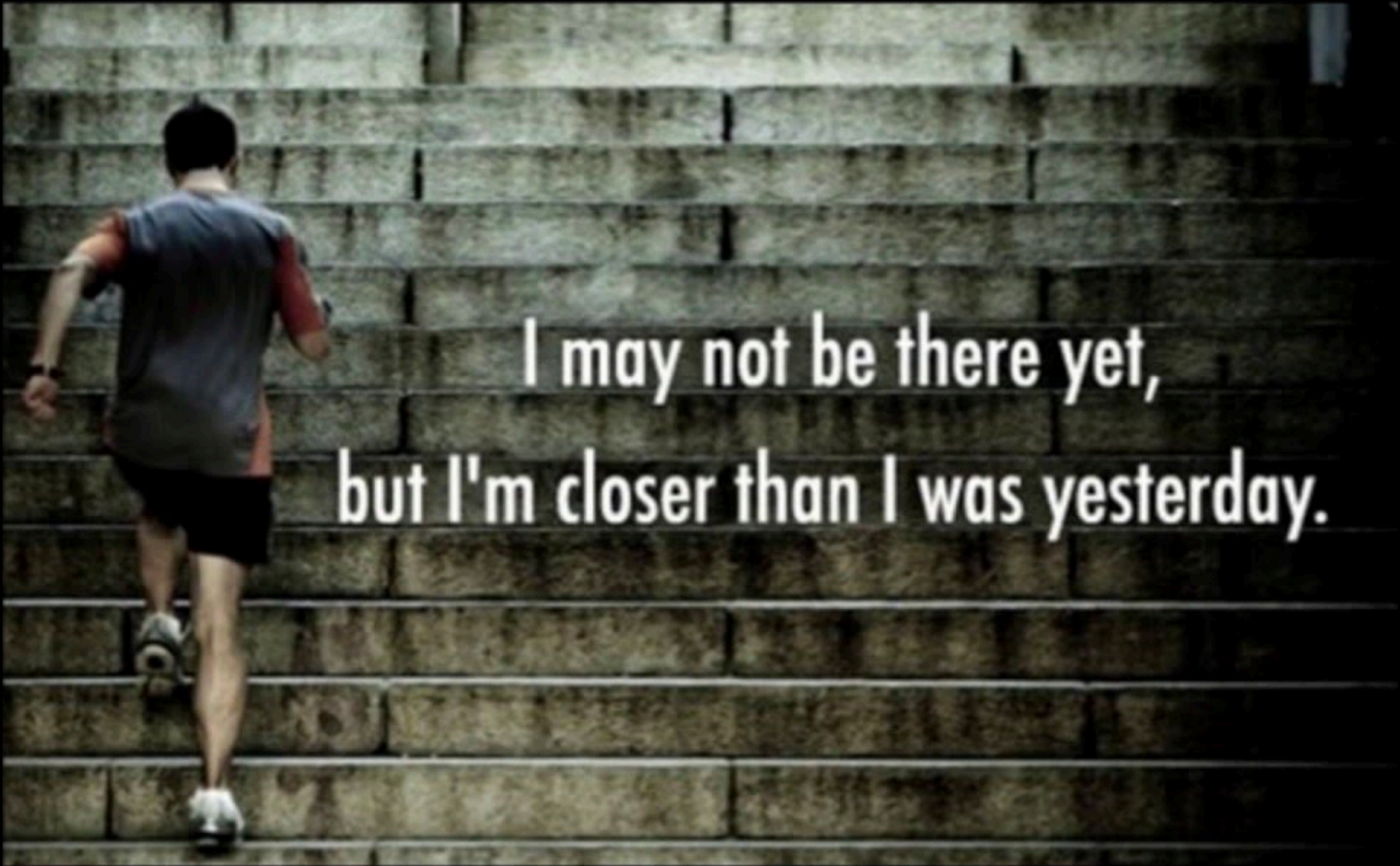with no loss of enthusiasm."

WINSTON CHURCHILL

# REDUNDANCY

I may not be there yet,
but I'm closer than I was yesterday.

# TED Slide Style

Don't use this for a technical talk

TED is great inspiration for **speaking** style
- but the slide format is mainly relevant for "motivational" talks


Similarly, much of the advice for making great slides online is not relevant!
- They are for a business audience!

# How many slides?

I typically aim for ~1 minute per slide

Varies depending on the depth of information per slide and whether you use "real" animation or multi-slide animation

# Know your audience!

Slide format will be very different…

Class

 - lots

Talk a

 - Pre

Pitchi

 - Foc

Talk a

 - The

P != NP

Each company / org will have its own "culture"

# P != NP

Blah blah

I proved it using the X conjecture

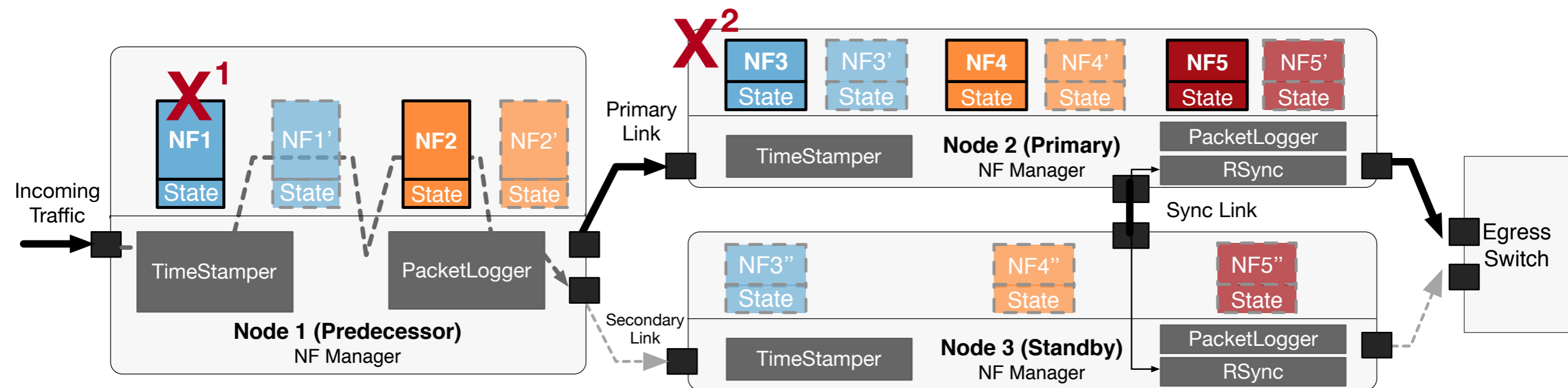Blah blah

# Diagrams and Graphs Tips

# Diagrams

Visual representations of your algorithm, system, or approach are always helpful

- Make the paper easier to understand
- Break up large chunks of text



Find a tool that works for you

- My lab used Omnigraffle (mac only), switching to **diagrams.net**
- Use something consistently so you become more efficient

# Bad Diagrams

Is this a good system diagram?

# Common Problems

Useless content

Bad color choices
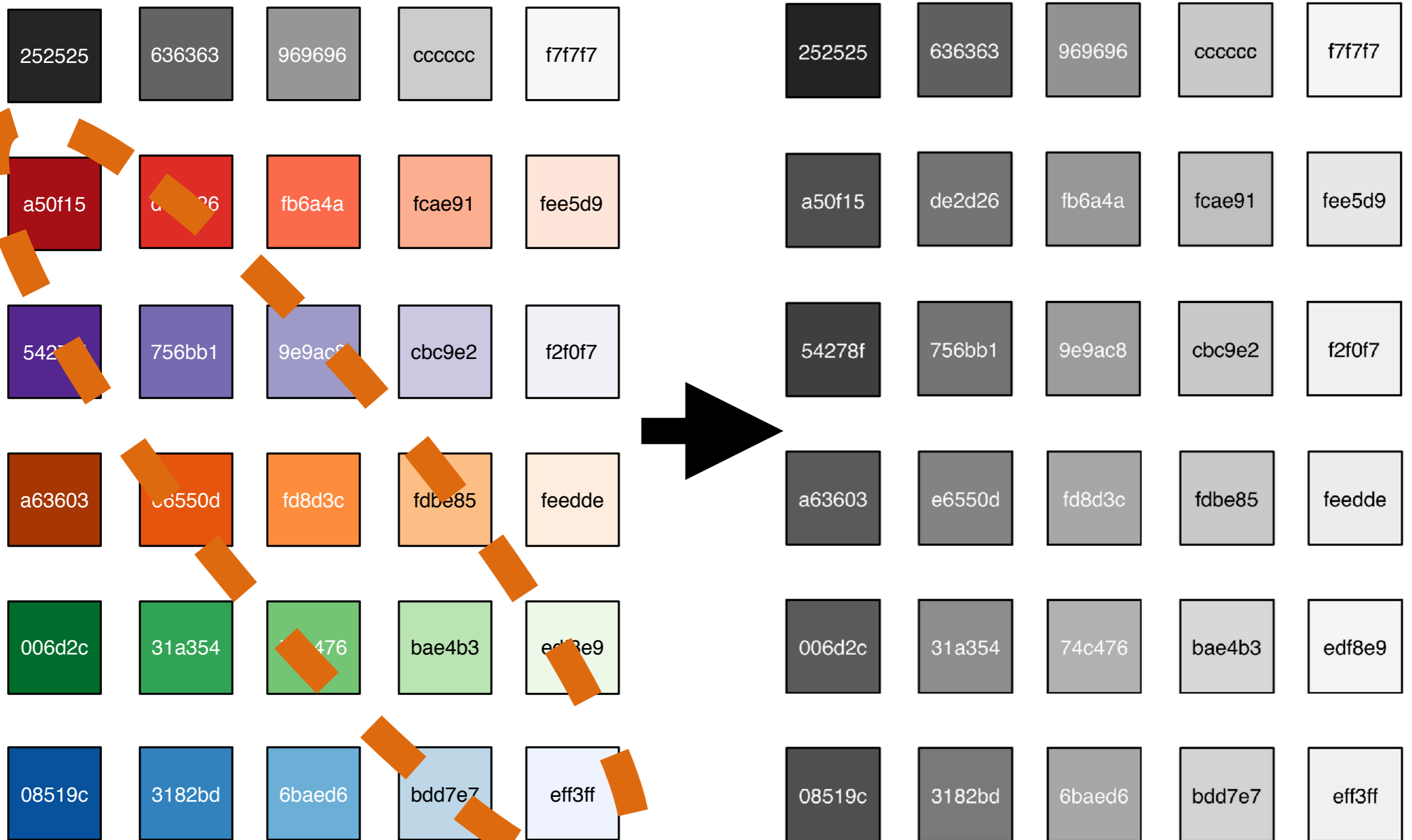- Indistinguishable, childish, etc

Fonts need to be bigger!

No caption to explain
- I prefer useful captions instead of "title" captions

# Color to BW

| | | | | |
|---|---|---|---|---|
| 252525 | 636363 | 969696 | cccccc | f7f7f7 |
| a50f15 | de2d26 | fb6a4a | fcae91 | fee5d9 |
| 54278f | 756bb1 | 9e9ac8 | cbc9e2 | f2f0f7 |
| a63603 | e6550d | fd8d3c | fdbe85 | feedde |
| 006d2c | 31a354 | 74c476 | bae4b3 | edf8e9 |
| 08519c | 3182bd | 6baed6 | bdd7e7 | eff3ff |

➡

| | | | | |
|---|---|---|---|---|
| 252525 | 636363 | 969696 | cccccc | f7f7f7 |
| a50f15 | de2d26 | fb6a4a | fcae91 | fee5d9 |
| 54278f | 756bb1 | 9e9ac8 | cbc9e2 | f2f0f7 |
| a63603 | e6550d | fd8d3c | fdbe85 | feedde |
| 006d2c | 31a354 | 74c476 | bae4b3 | edf8e9 |
| 08519c | 3182bd | 6baed6 | bdd7e7 | eff3ff |

# Color to BW

| | | | | |
|---|---|---|---|---|
| 252525 | 636363 | 969696 | cccccc | f7f7f7 |
| a50f15 | de2d26 | fb6a4a | fcae91 | fee5d9 |
| 54278f | 756bb1 | 9e9ac8 | cbc9e2 | f2f0f7 |
| a63603 | e6550d | fd8d3c | fdbe85 | feedde |
| 006d2c | 31a354 | 74c476 | bae4b3 | edf8e9 |
| 08519c | 3182bd | 6baed6 | bdd7e7 | eff3ff |

➡️

| | | | | |
|---|---|---|---|---|
| 252525 | 636363 | 969696 | cccccc | f7f7f7 |
| a50f15 | de2d26 | fb6a4a | fcae91 | fee5d9 |
| 54278f | 756bb1 | 9e9ac8 | cbc9e2 | f2f0f7 |
| a63603 | e6550d | fd8d3c | fdbe85 | feedde |
| 006d2c | 31a354 | 74c476 | bae4b3 | edf8e9 |
| 08519c | 3182bd | 6baed6 | bdd7e7 | eff3ff |

# Relationships

Most good diagrams focus on the relationships between different components/concepts

A good diagram should help a reader use the "black box" technique to filter out aspects of the work that are less important to them
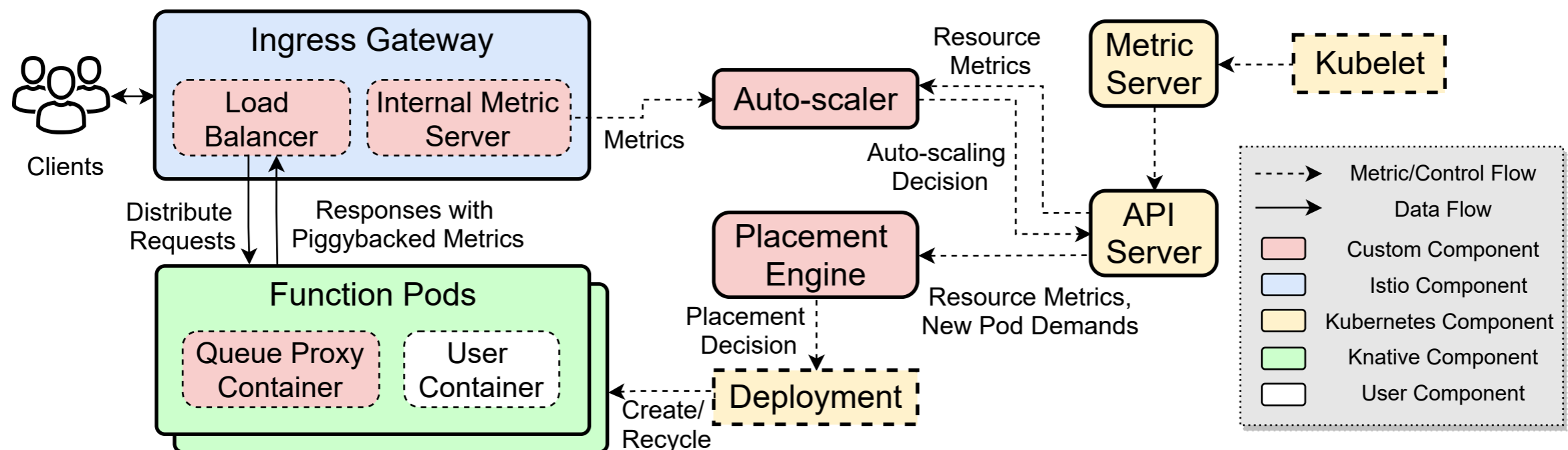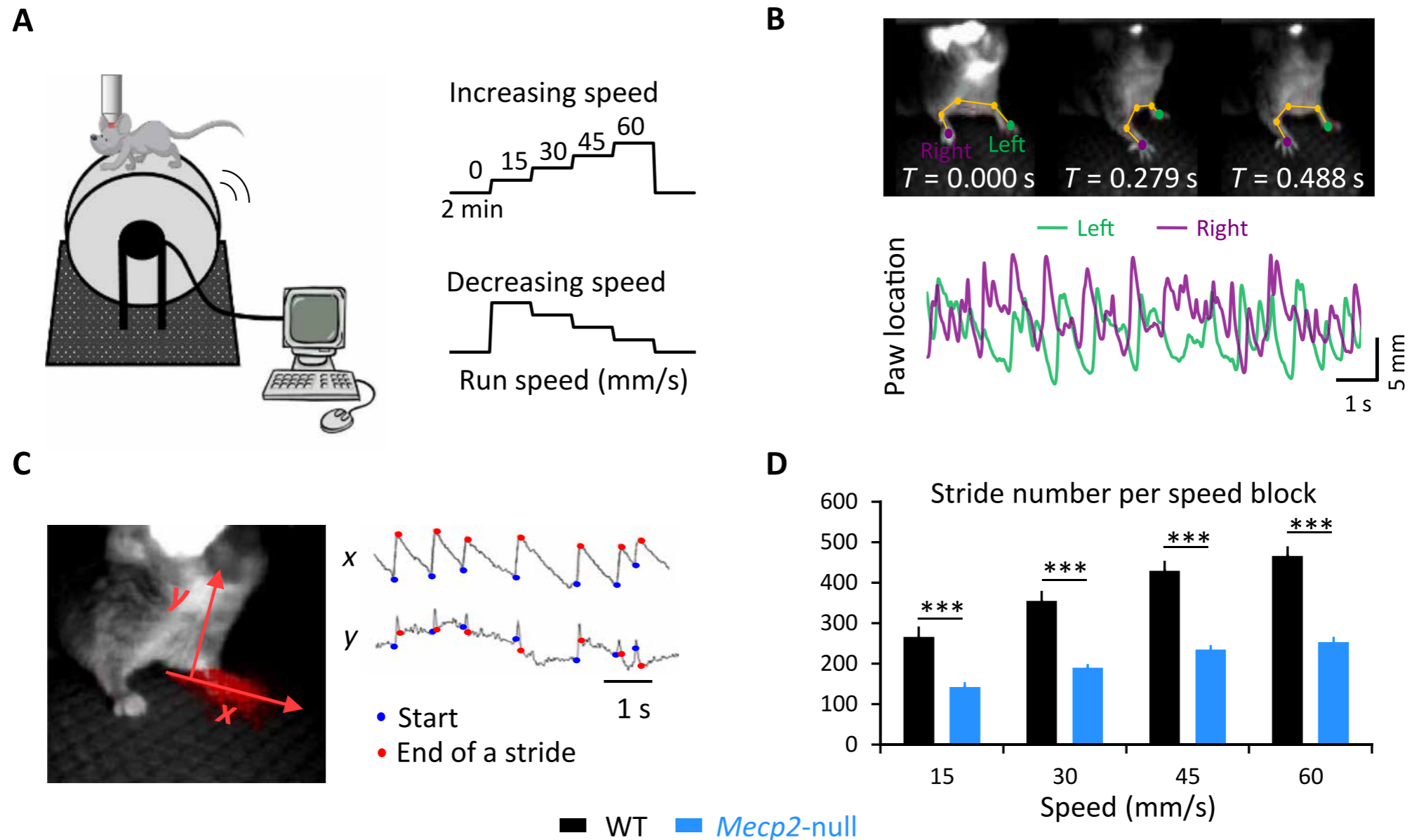


**Figure 1: Mu Overview.**

# Diagrams vs Results

## Diagrams can explain experimental setup



From "Motor training improves coordination and anxiety in symptomatic *Mecp2*-null mice despite impaired functional connectivity within the motor circuit" by Simha and friends

| Query type | Scheme (References) | Approach | # of parties | Threats | | $S$ leakage | | Scale | | | Crypto | | | Network | | Unique feature |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | Adversarial $Q$ | Adversarial $S$ | Init | Query | Updatable? | Implemented? | Scale tested | Crypto type | Insert: # ops | Query: # ops | # round trips | Data sent | |
| Equality | Arx-EQ [14] | Legacy | 2 | — | ◐ | ○ | ◕ | ● | ✔ | ◐ | ● | ● | ● | ● | ● | legacy compliant |
| | Kamara-Papamanthou [106] | Custom | 2 | — | ◐ | ○ | ◕ | ● | — | — | ● | ○ | ○ | ● | ● | parallelizable |
| | Blind Storage [100] | Custom | 2 | — | ◐ | ○ | ◕ | ● | ✔ | ◐ | ● | ◐ | ● | ◐ | ● | low $S$ work |
| | Sophos ($\Sigma o\phi o\varsigma$) [101] | Custom | 2 | — | ◐ | ● | ◕ | ● | ✔ | ◐ | ○ | ◐ | ● | ● | ● | **Refresh** w/ **Insert** |
| | Stefanov et al [107] | Custom | 2 | — | ◐ | ○ | ◕ | ● | ✔ | ◐ | ○ | ○ | ○ | ● | ● | **Refresh** w/ **Insert** |
| | vORAM+HIRB [120] | Obliv | 2 | — | ◐ | ○ | ○ | ○ | ✔ | ● | ● | ○ | ○ | ○ | ◕ | history independ. |
| | TWORAM [121] | Obliv | 2 | — | ◐ | ○ | ○ | — | — | — | ◐ | ○ | ○ | ◐ | ◕ | const round |
| | 3PC-ORAM [124] | Obliv | 3 | ◐ | ◐ | ○ | ○ | ● | ✔ | ◕ | ● | ○ | ○ | ○ | ◕ | dual $S$ |
| Boolean | DET [15], [92] | Legacy | 2 | — | ◐ | ◑ | ◑ | ● | ✔ | ● | ● | ◐ | ◐ | ● | ● | supports JOINs |
| | BLIND SEER [16], [17] | Custom | 3 | ● | ● | ○ | ◐ | ◐ | ✔ | ● | ◐ | ◐ | ◐ | ○ | ◕ | hide field, $r_i$'s |
| | OSPIR-OXT [18]–[21], [104] | Custom | 3 | ● | ◐ | ○ | ◐ | ● | ✔ | ● | ◐ | ◐ | ◐ | ◑ | ● | excels w/ small $r_1$ |
| | Kamara-Moataz [102] | Custom | 2 | — | ◐ | ○ | ◐ | ○ | — | — | ◐ | ◐ | ○ | ● | ◕ | relational SPC |
| Range | OPE [93]–[95] | Legacy | 2 | — | ◐ | ● | ● | ● | ✔ | ● | ● | ● | ● | ● | ● | leak some content |
| | Mutable OPE [97] | Legacy | 2 | — | ◐ | ● | ● | ● | ✔ | ● | ● | ○ | ○ | ○ | ◐ | interactive |
| | Partial OPE [111] | Custom | 2 | — | ◐ | ○ | ● | ● | ✔ | ◐ | ● | ● | ● | ◐ | ● | fast insertions |
| | Arx-RANGE [110] | Custom | 2 | — | ◐ | ○ | ◐ | ● | ✔ | ◐ | ◐ | ○ | ◕ | ● | ○ | non-interactive |
| | SisoSPIR [22] | Obliv | 3 | ◐ | ◐ | ○ | ○ | ○ | ✔ | ● | ● | ● | ● | ○ | ◕ | split, non-colluding $S$ |
| Other | GraphEnc$_1$ [116] | Custom | 2 | — | ◐ | ◐ | ◕ | ◐ | ✔ | ◐ | ● | ● | ● | ● | ◕ | approx. graph dist. |
| | GraphEnc$_3$ [116] | Custom | 2 | — | ◐ | ◐ | ◐ | ◐ | ✔ | ◐ | ○ | ● | ● | ● | ● | approx. graph dist. |
| | Chase-Shen [109], [126] | Custom | 2 | — | ● | ○ | ◐ | ○ | ✔ | ◕ | ● | ● | ● | ◐ | ● | substring search |
| | Moataz-Blass [123] | Obliv | 2 | — | ◐ | ○ | ○ | ● | ✔ | ● | ● | ○ | ○ | ○ | ◕ | substring search |

TABLE II

SUMMARY OF THE SECURITY, PERFORMANCE, AND USABILITY OF BASE QUERIES. $Q$ AND $S$ DENOTE THE QUERIER AND THE SERVER, RESPECTIVELY. WE PRESUME THAT THE ADVERSARY KNOWS THE DATABASE SIZE $d$ AND THE LENGTH OF EACH RECORD. FOR SYSTEMS THAT EITHER DO NOT SUPPORT INSERT OR USE A SIDE INDEX, THE INSERT COST IS THE AMORTIZED COST OF ADDING A SINGLE RECORD DURING **Init**. LEGENDS FOR EACH COLUMN FOLLOW. IN ALL COLUMNS EXCEPT "INIT/QUERY LEAKAGE," BUBBLES THAT ARE MORE FILLED IN REPRESENT PROPERTIES THAT ARE BETTER FOR THE SCHEME.

SCALE TESTED
● – BILLIONS
◐ – MILLIONS
◕ – THOUSANDS

UPDATABLE
● – INSERT IN MAIN INDEX
◐ – BUILD SIDE INDEX
○ – NOT SUPPORTED

THREATS
● – MALICIOUS
◐ – SEMI-HONEST

DATA SENT
(BEYOND RESULTS)
● – CONSTANT
◐ – ADDITIVE POLYLOG($d$)
◕ – MULT. POLYLOG($d$)
○ – EVEN MORE

INIT/QUERY LEAKAGE
(SEE SECTION II-E)
● – ORDER/CONTENTS
◑ – EQUALITY
◐ – PREDICATE
◕ – IDENTIFIER
○ – STRUCTURE

TYPE OF CRYPTO
● – SYMMETRIC
◐ – BATCHED OR PRE-COMPUTED PUBLIC-KEY
○ – PUBLIC-KEY

CRYPTO OPS PER RECORD
● – CONSTANT
◐ – # KEYWORDS
○ – LOGARITHMIC

ROUND TRIPS
● – 1
◑ – 2
◐ – CONSTANT
○ – LOGARITHMIC

From "Sok: Cryptographically protected database search" by Yerukhimovich and friends

# Condensing Information

A diagram is a way of extracting out the key concepts of a complex system or phenomenon

If you don't know how to draw it, you probably don't fully understand it yourself!

Practice sketching out your systems/algorithms/concepts as diagrams to help in conversations with your advisor/collaborators

- Communication is much more efficient with a whiteboard!
- Don't let your advisor do all the drawing!

# Let's make a diagram…

Goals:

Clean, consistent shapes

Clear connections between components

Useful colors

Looks good in print

# Running Experiments

*I ran an experiment where I evaluated the effectiveness of my new operating systems, TimOS*

*When running BenchmarkX, TimOS had a score of 250,000!*

Am I done? What do you need to know to interpret these results?

# An experiment is a sample

The result of an experiment is not "truth"

It is a **sample**, drawn from a **population** of possible results for your combination of Task, System, and Environment

- In some cases one experiment gives you many samples, e.g., accuracies of 10,000 images classified, response time of 10M requests, etc

# Data Analysis

You need to look closely at your data!

## Is it consistent?
- Always repeat experiment where possible

## Does it match your expectation?
- You should have had a mental model for the expected result

# Detailed Statistics

Average is not enough

Often you need to understand the **distribution**

# Bar vs Line Graphs

Bar graphs let you break the "axis must be consistently spaced" rule

Or use bars when measurements along x-axis are not ordered (categories of values, experiment #, etc)

Whenever you connect the dots, think about their meaning

# More Bars

## For % of whole, why not use pie chart?



Figure 7: Comfort with monitoring types (**Q28**).

# Histogram or PDF

How many samples occurred within each bucket range?

Or, what is the probability of hitting a certain range?
- Normalized histogram or continuous PDF



session_duration_seconds

# CDF

Cumulative Distribution Function

What percent of samples had this score or lower

## How to read?

- Find quartiles and 99 percentile on y-axis and go right until you hit the line, then down to the x-axis



Legend:
Case 1
Case 2
SPIKE (1 SmartNIC)
SPIKE (2 SmartNICs)
SLA threshold

% of requests

Percentage of response time distribution
**response time (ms)**

Concurrency
RPS
Mu

Avg:
7765
2605
1020

# When to use CDF/PDF

Use when outliers/tail are important

Histogram is often easier to read

But, comparing multiple histograms takes a lot of space

CDF is useful when comparing multiple distributions in one plot



Avg:
588
880
952

# Box-Plots

Quickly display summary statistics and allow easy comparison



Whiskers: Min & Max
center line: Median (50%ile)
top/bottom: 75/25th %ile

# Boxplot Variations

Whiskers aren't always max/min

Might show another statistics such as 2nd and 98th percentile, with outlier points explicitly shown

Paper should specify what the whiskers show

# Violin Plots

Like a box plot, but shows the full distribution

More descriptive

Useful if comparisons have a fundamentally different distribution

# ???

# Same Stats, Different Graphs

https://www.autodesk.com/research/publications/
same-stats-different-graphs

# Graphs

Stick with "standard" graph types

- Unless you have a good reason
- Standards: Scatter plot, line plot, bar plot
- Avoid pie charts and "infographic" look

Include error bars!

- Use standard deviation
  or confidence interval

**Bar Chart**

**Line Graph**

Every graph in this paper was a radar chart… odd

# Graphs

What tool do you use to make plots?
- gnuplot, matplotlib, seaborn, tableau, matlab, R…

Avoid tools like Excel
- Most papers I review with Excel graphs I reject (not usually because of the graphs, but it is a sign of amateur-ness)
- OK for initial data exploration by yourself

Use the same tool as your lab-mates
- Have lab scripts for making beautiful graphs

# Jupyter Notebooks

This is the "right" way to do it

Browser-based (or VS Code extension!) editor and python runtime

Parse, analyze, and visualize data all in one place

Easy to share with others

- Github displays rendered notebooks
- Google Collab allows shared notebooks / execution environments

# Scripting Experiments

More experiments is more data

More data is more information

More information is more value
- As long as you have the time to process it

You should run as many experiments as possible
- Uses scripts to automate gathering AND analyzing data as much as possible
- Try to make general purposes scripts you can reuse in the future

You will need to repeat experiments, so make it easy!

# Basic Shell Scripting

Often output of experiments is messy

Bash scripts that combine basic Unix utilities are usually the fastest way to pre-process your data

Or you can do this all in Python, but may be simpler as part of your experiment running scripts

*grep "exp1" file.data | awk '{print $3}' | sort -n | tail*

sed - replace

awk - processing rows / columns in a file

grep - filtering

# What is wrong with this?
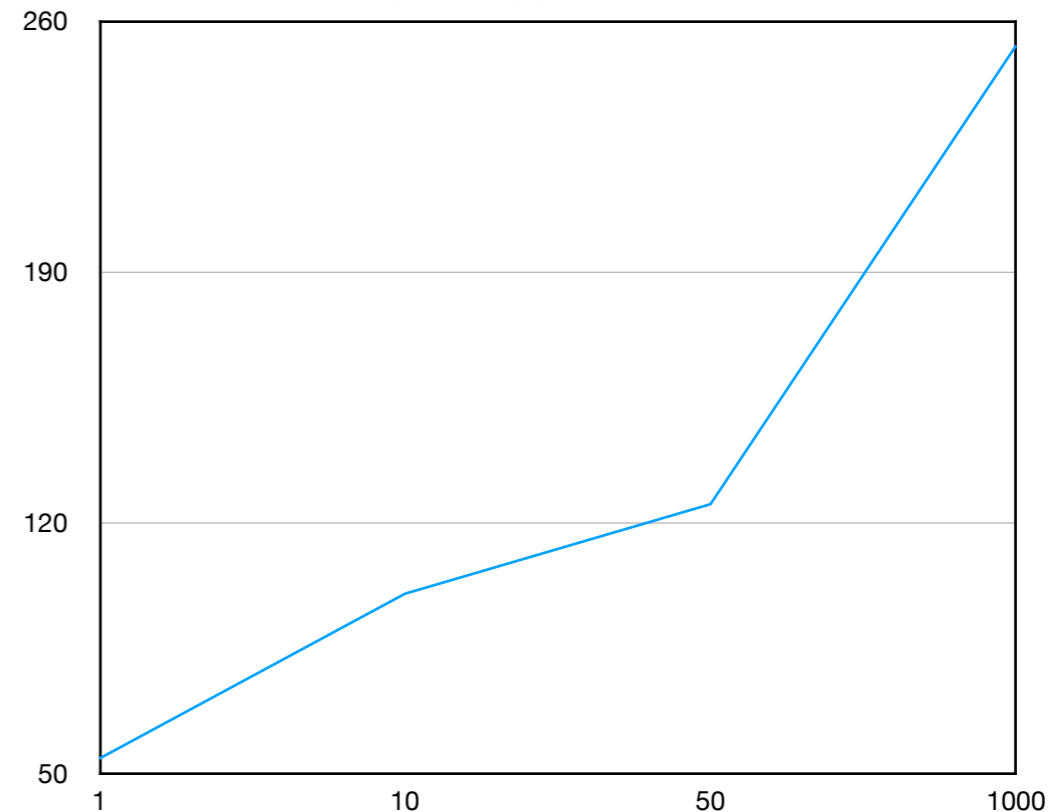
# What is wrong with this?

Needs axis labels

No error bars

Irregular x-axis

Non-zero y-axis

Font sizes too small



Aspect ratio may be awkward in paper layout
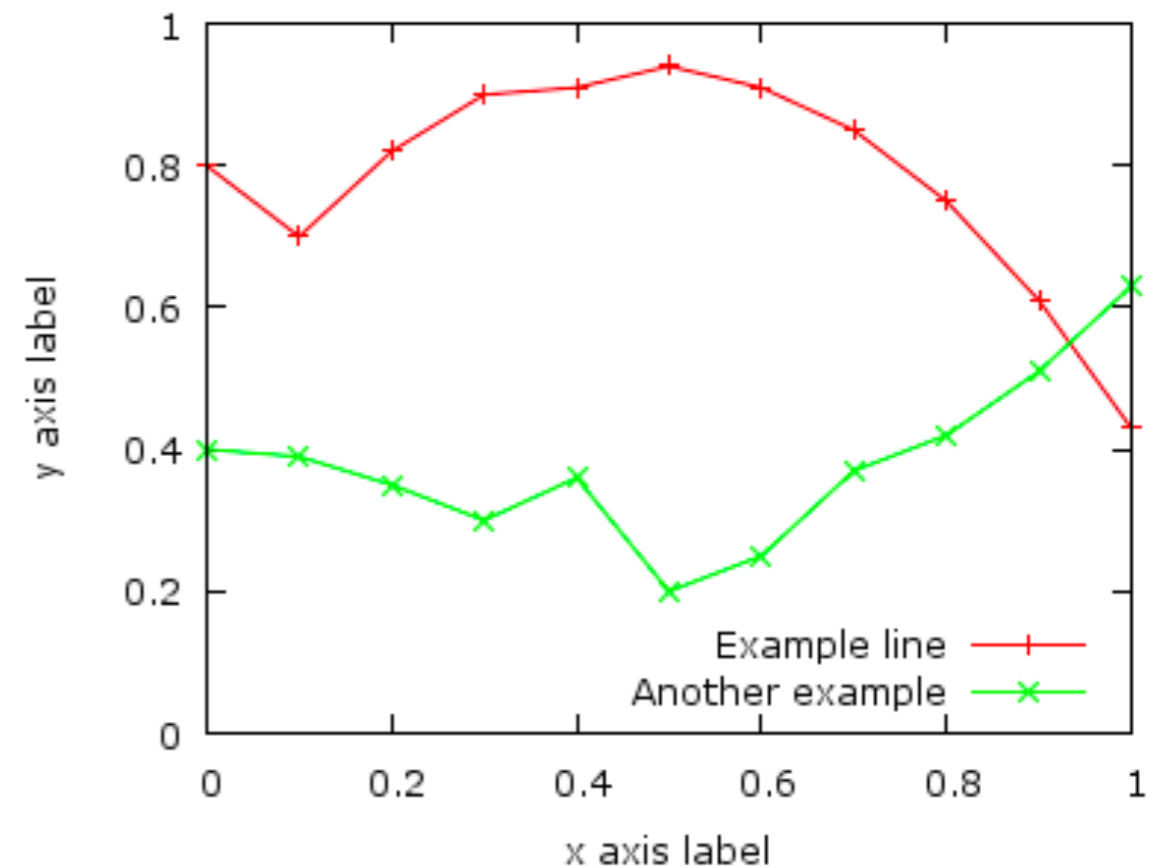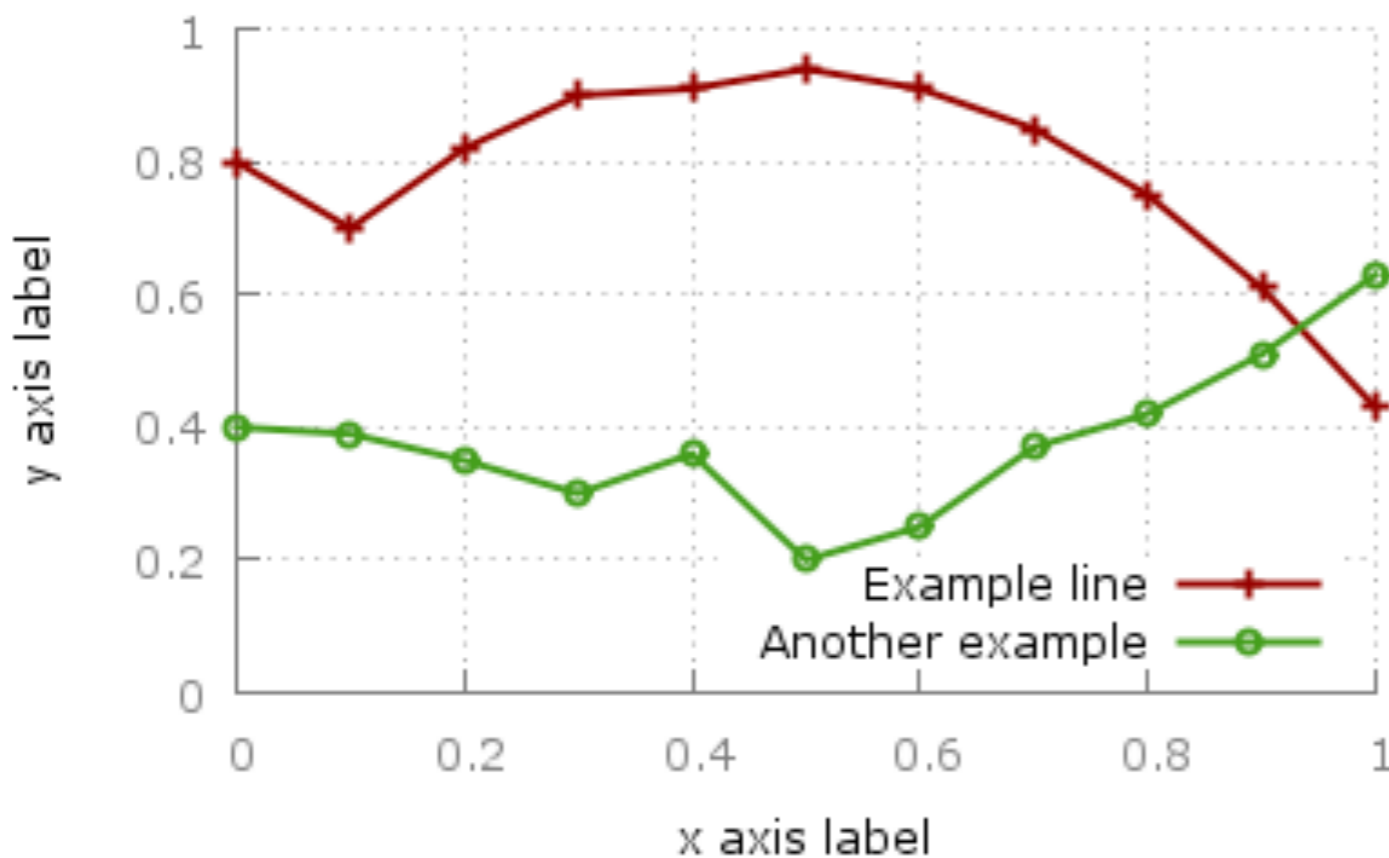- Typically either want 2 square images or 1 wide image for 2-column layout

# Graph Tips

## Use:

- Thick lines
- Very large fonts
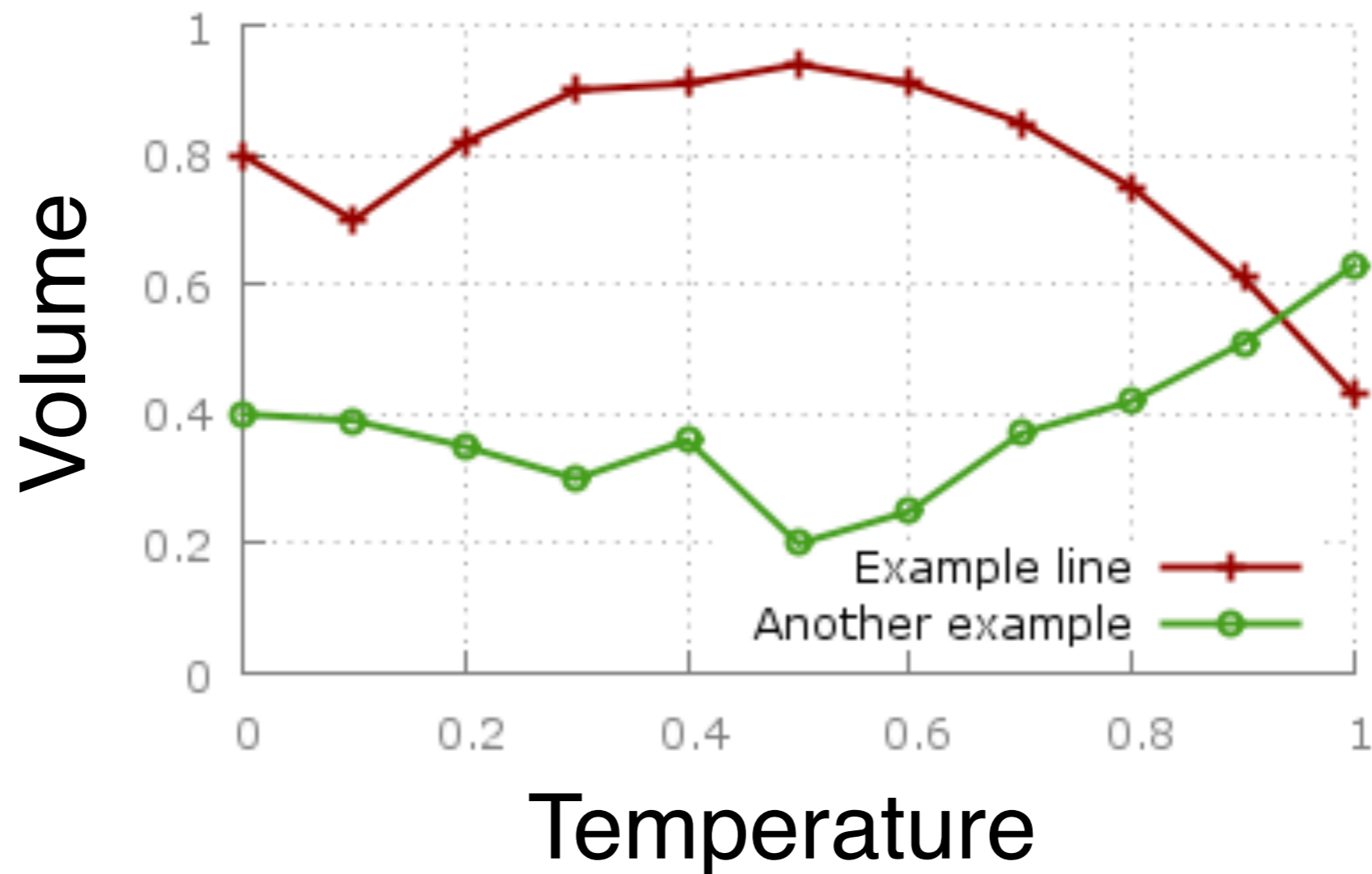- Axis labels
- Wide format

## Avoid:

- Similar colors (check BW!)
- Non-0 starting axes
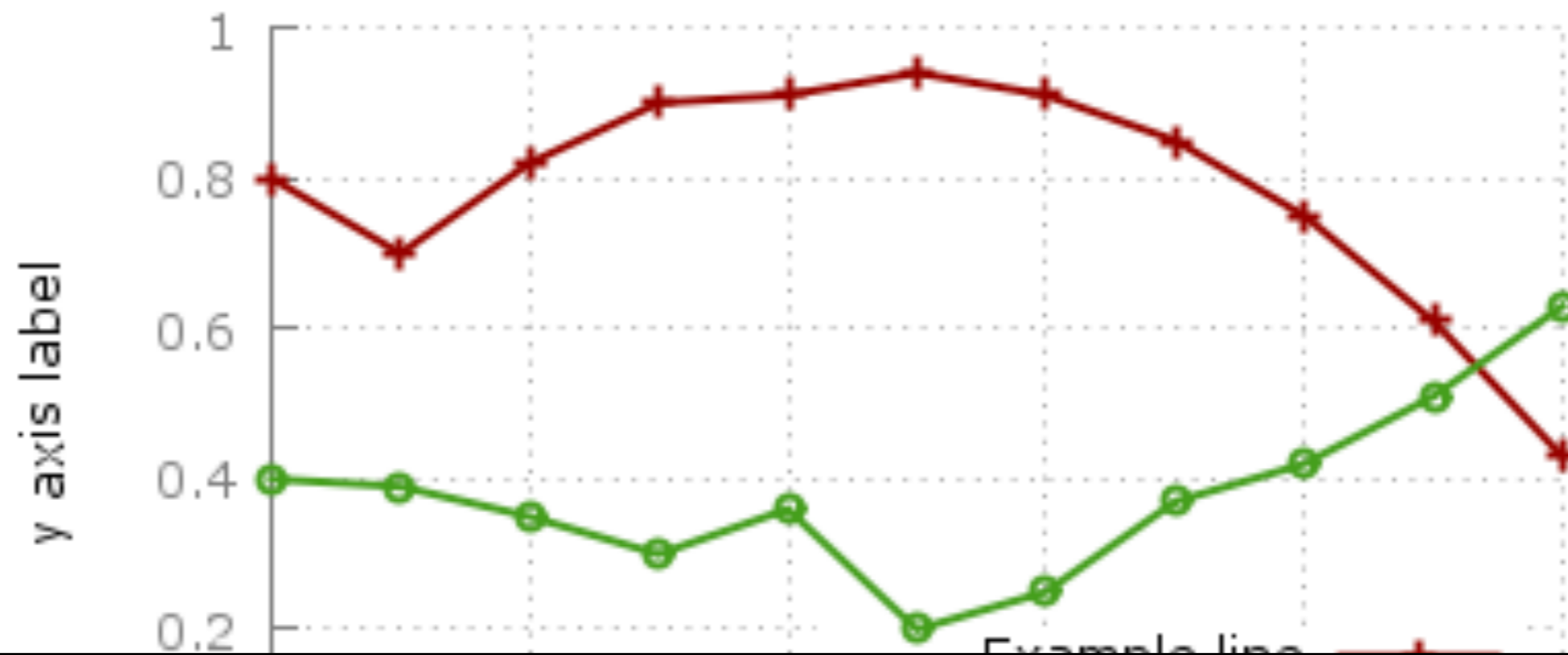- Titles (if you have caption)
- Square aspect ratio

# Line Graphs…

My high school science teacher would (correctly) fail me for making this graph… why?
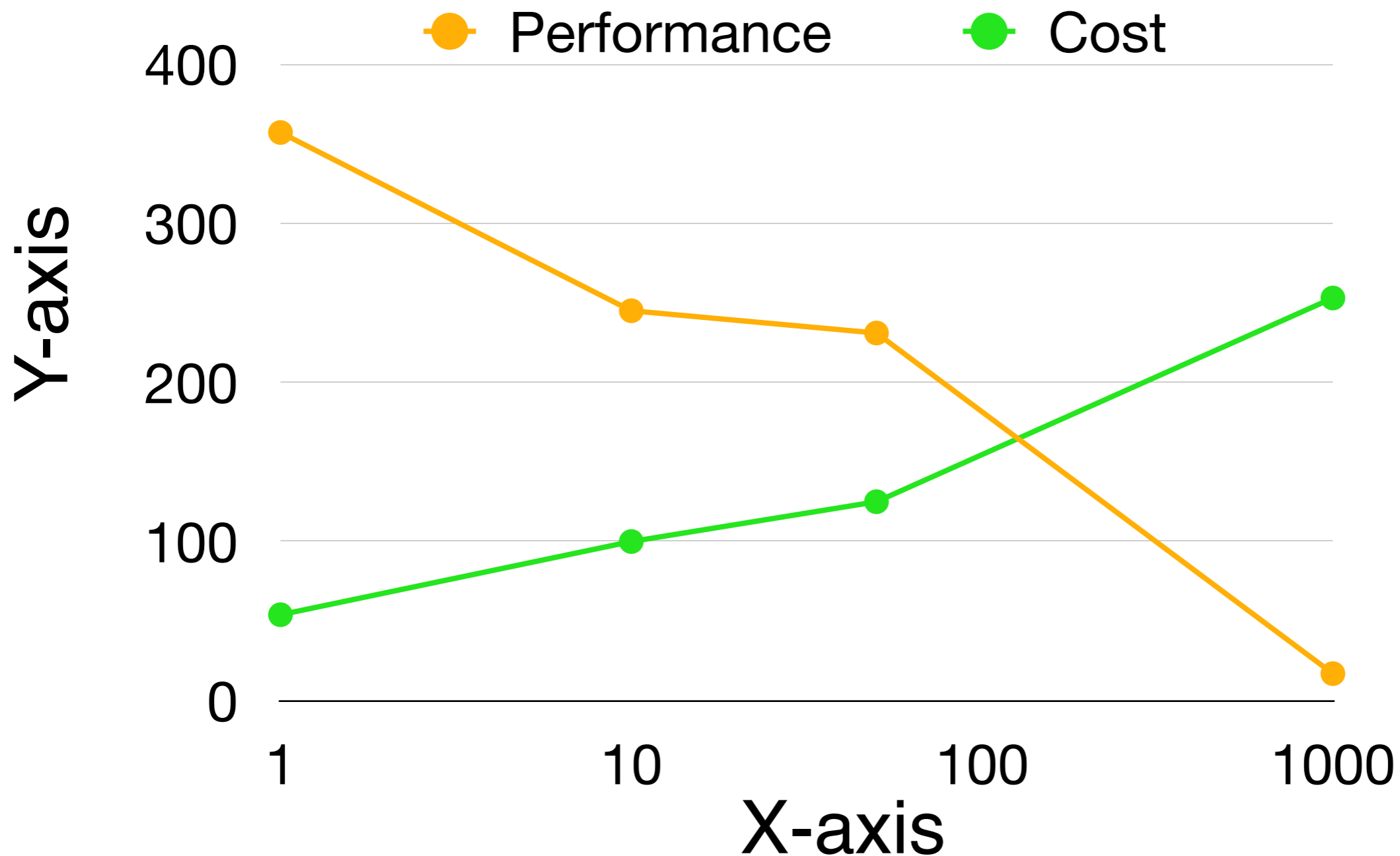
# Line Graphs…

My high school science teacher would (correctly) fail me for making this graph… why?



"Connect the dots" is really unscientific!
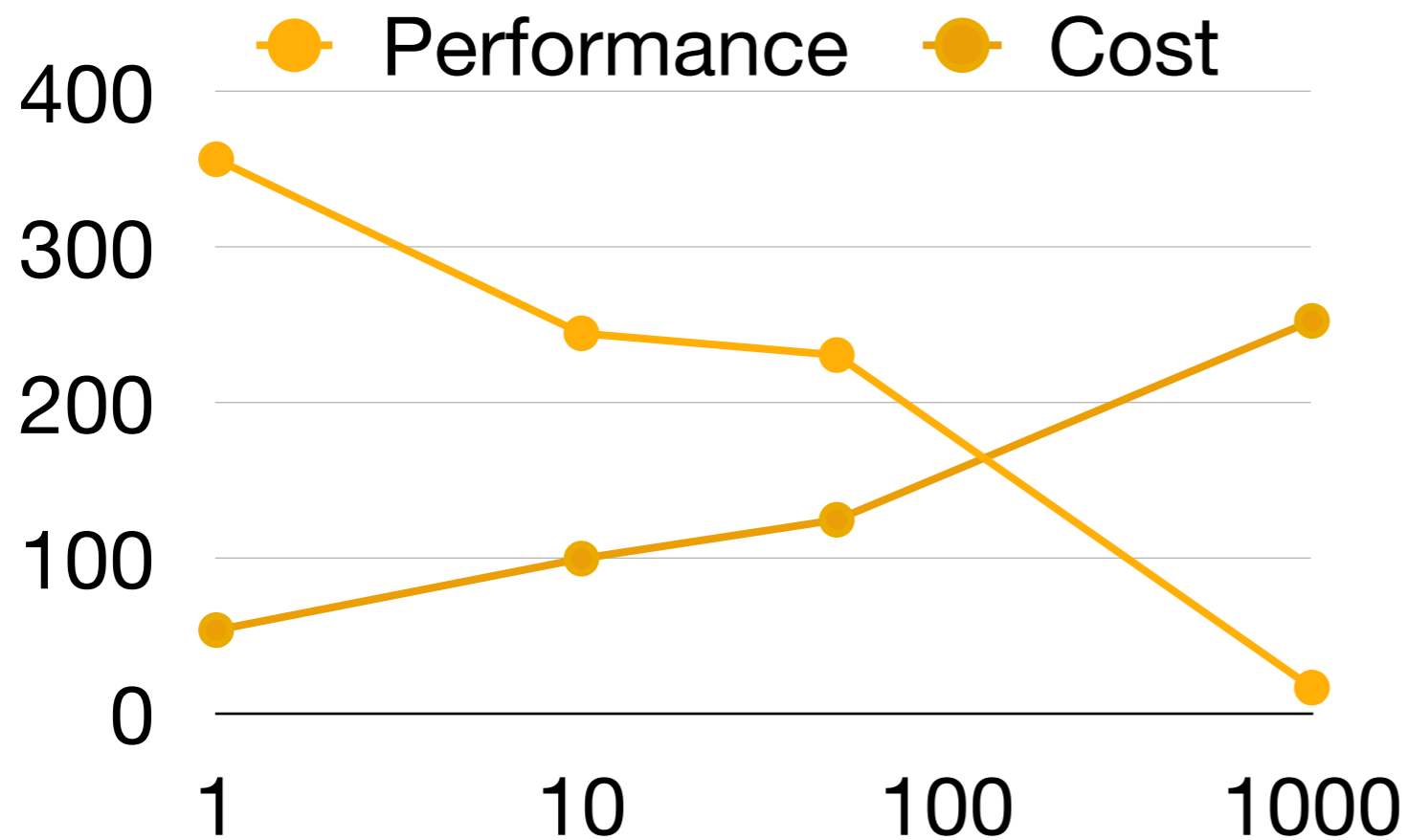Trend lines are way better!

But this is what the community expects!

# Problem with this?

# Accessibility

## Color blindness
- 1 in 12 men, 1 in 200 women



THE AMERICAN FLAG
AS SEEN
I.-BY MOST PEOPLE
II.-BY RED-BLIND PERSONS
III.-BY GREEN-BLIND PERSONS
IV.-BY VIOLET-BLIND PERSONS
V.-BY TOTALLY COLOR-BLIND PERSONS

# Accessibility

## Color blindness

- 1 in 12 men, 1 in 200 women

## Blindness

- about 8 million people in the US have a visual disability

Many types of disabilities to be aware of: sight, sound, touch, mobility

Remember these!

THE AMERICAN FLAG
AS SEEN
I.- BY MOST PEOPLE
II.- BY RED-BLIND PERSONS
III.- BY GREEN-BLIND PERSONS
IV.- BY VIOLET-BLIND PERSONS
V.- BY TOTALLY COLOR-BLIND PERSONS

# Visual Recipes

Slides:
- Mixture of text and images
- Keep bullets simple, fonts clear
- Use animation sparingly for emphasis

Diagrams:
- Pick the right level of abstraction; focus on relationships
- Use LARGE fonts!
- Be sure colors work in B&W

Graphs:
- Use easy to understand plot types
- Use thick lines and be sure they are distinguishable
- Use LARGE fonts!
- Be sure colors work in B&W