# Building your own Kubernetes and Docker

Erik Bernhardsson

# I'm Erik Bernhardsson

- Founder of Modal Labs
- Built the music recommendation system at Spotify
- I tweet sometimes: @bernhardsson
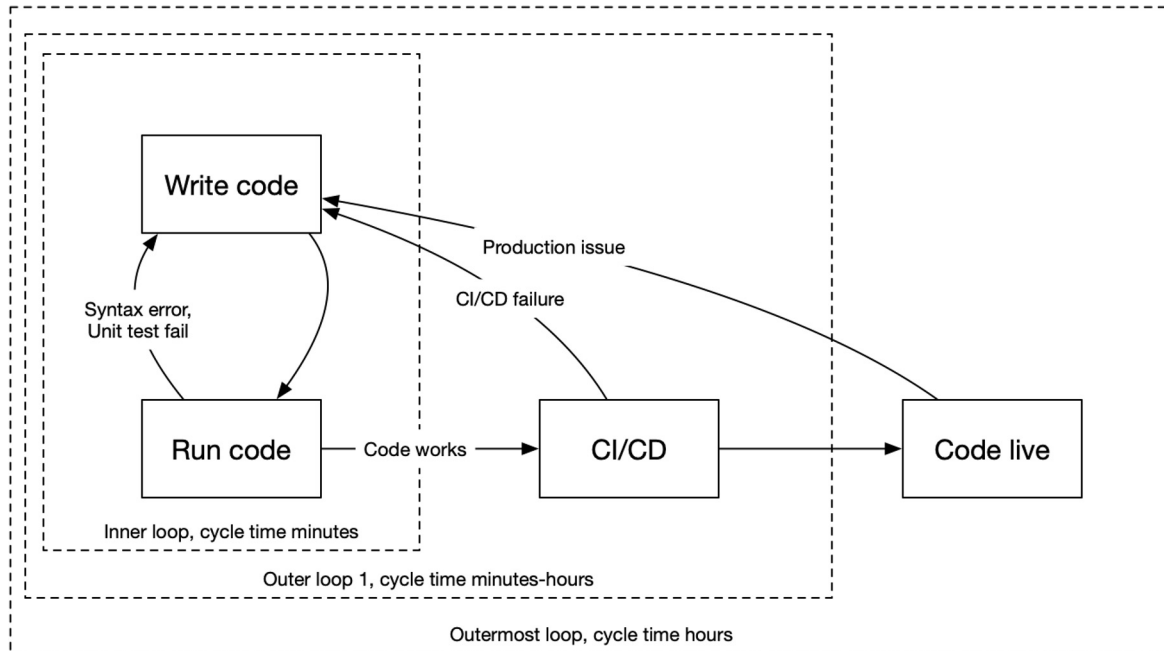- I blog very occasionally: https://erikbern.com

# All just wanted to make data teams more productive!

- How to productionize jobs
- How to scale things out
- Scheduling things
- How to use GPUs and other hardware

# What do I mean by eng productivity

A set of nested for-loops of writing code

# Frontend

- Write code in one window
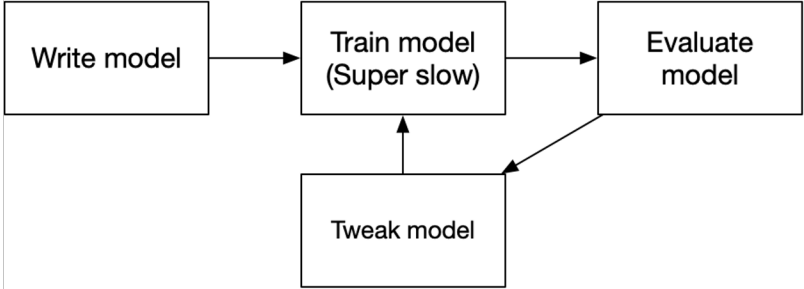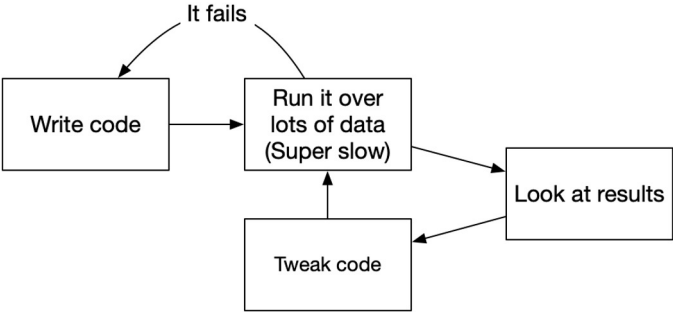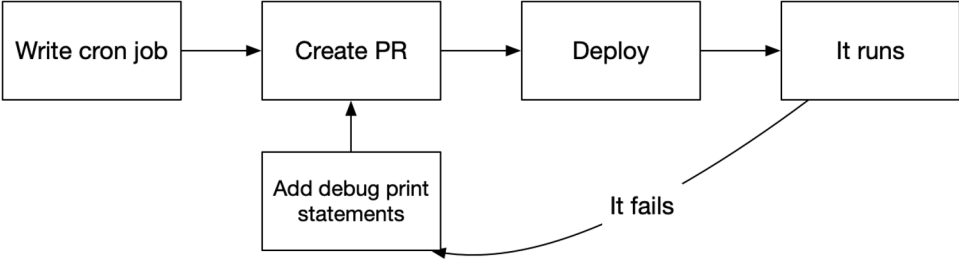- Look at the website in another window

# Backend

1. Write code
2. Does it compile?
3. Does it pass unit tests?
4. Ship it

# Data has super long feedback loops

# Let's put infrastructure into the feedback loop

If we get most of this to happen in the cloud instead:

- Moves a lot of stuff from an outer loop into an inner loop
- If env is always the same, it reduces a whole set of things that can break
- We have infinite compute power and storage
- Never have to think about drivers and GPUs

# What are containers?

- Represent all dependencies as a Linux root filesystem
- Have a bunch of stuff for resource management (and to a limited extent, security)

# Cracking open a Docker container

```
$ docker pull python

$ docker run -d python
sleep infinity

$ docker export
b0aa33209370 > python.tar

$ tar tvf python.tar
```
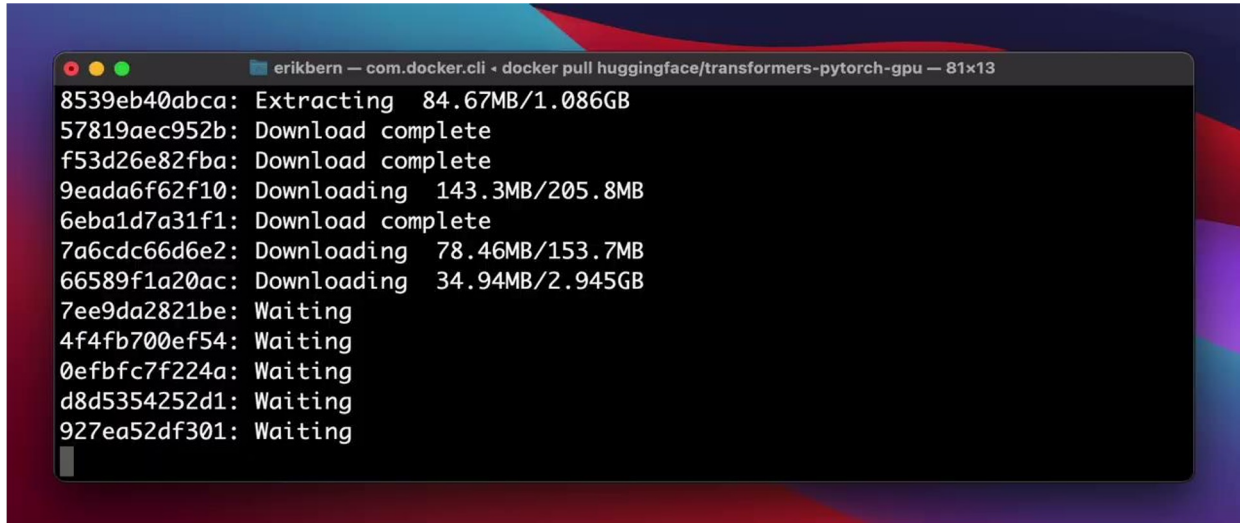
# How to launch a container on a remote host

1. Pull down an image: a few sec to a few minutes
2. Start the image: a couple of seconds

# The average container image has a lot of junk

Eg the `python` container from Dockerhub:

- 870MB large
- 29,772 files
  - /usr/share/locale: 1,553 files
  - /usr/share/doc: 3,210 files
  - /usr/share/perl: 1,389 files
  - /usr/share/man: 3,050 files

# What do we actually need to run something?

`$ python3 -c 'import sklearn'`

- 3,043 calls to `stat`
- 1,073 calls to `openat`

Lots of file system operations!

But only a small number of unique files are accessed.

# It would be nice to avoid Docker

`runc` is a nice utility:

- Point it at a root file system
- It runs a container!
- Not absurdly complex (~50k lines of Go)

# Basic container runner that avoids docker pull:

After building the image: **docker save** to a network drive

When running the container: **runc** with a root filesystem over the network



© Vernier Software & Technology

# This is still pretty slow though!

- Python does thousands of file system operations sequentially
- NFS latency is a few milliseconds!

This adds up to like 10 seconds!

If we want to do this in seconds, we have a fraction of a millisecond for each file system operation.

Rough latency numbers:

- S3: 10-20ms
- NFS: 1-2ms
- EBS: 0.5-1ms
- SSD: 100-200 µs

# Can we cache things locally?

- SSD latency: ~100 µs (0.1ms)
- Same image: almost the same files are read every time
- Different image: still almost the same files every time!

# Unrelated images have a lot of overlap!

pytorch/pytorch

18,240
unique files

6,415

78,287
unique files

huggingface/transformers-pytorch-gpu

8,716

10,695

2,228

smy20011/dreambooth

25,438
unique files

# How to cache efficiently: content-addressing

Indexes of filesystems

Storage

```
/etc/xyz  0abc730f
/lib/x86_64-linux-gnu/libc-2.31.so  67dd99b
/usr/share/python/os.py  388cb9d
/root/my_funny_script.py  6593bb9
```

```
/etc/qqq  bab933d
/lib/x86_64-linux-gnu/libc-2.32.so  837733c
/usr/share/python/os.py  388cb9d
/root/hello.txt  b99d9b9
```

```
/etc/qqq  bab933d
/lib/x86_64-linux-gnu/libc-2.31.so  67dd99b
/usr/share/python/os.py  544c453
/root/my_funny_script.py  6593bb9
```

```
/0a/bc/730
/38/8c/b9d
/54/4c/453
/65/93/bb9
/67/dd/99b
/83/77/33c
/b9/9d/9b9
/ba/b9/33d
```

# How do we make this work with containers?

Build our own file system:

- Not super hard with FUSE!
- You can even do it in Python
- A lot easier if the file system is read-only

# FUSE operations we need to implement

`open`

`read`

`release`

`readdir`

`readdirplus`

# Handle the indirection when reading files

Keep an index in memory that maps file system paths to

1. The hash of the content
2. A **struct stat**

```
struct stat {
    dev_t     st_dev;       /* ID of device containing file */
    ino_t     st_ino;       /* Inode number */
    mode_t    st_mode;      /* File type and mode */
    nlink_t   st_nlink;     /* Number of hard links */
    uid_t     st_uid;       /* User ID of owner */
    gid_t     st_gid;       /* Group ID of owner */
    dev_t     st_rdev;      /* Device ID (if special file) */
    off_t     st_size;      /* Total size, in bytes */
    blksize_t st_blksize;   /* Block size for filesystem I/O */
    blkcnt_t  st_blocks;    /* Number of 512B blocks allocated */
};
```

# When reading a file

1. Look up its hash in the index
2. See if it exists on local disk
   a. If not, fetch it, return its content, and store the file on local disk
   b. If it does exist, just return it

Index

Storage

/etc/xyz **0abc730f**
/lib/x86_64-linux-gnu/libc-2.31.so **67dd99b**
/usr/share/python/os.py **388cb9d**
/root/my_funny_script.py **6593bb9**

/0a/bc/730
/38/8c/b9d
/54/4c/453
/65/93/bb9
/67/dd/99b
/83/77/33c
/b9/9d/9b9
/ba/b9/33d

# Ok but how do we get the images into this?

We already build the containers in the cloud so that's a good starting point!

Super janky idea:

- Build images using Docker
- Then `docker save` to a temporary directory
- Then checksum of every file
  - Upload any file to NFS that we didn't have already
- Then build an index of path → (checksum, struct stat)
- Store the index on NFS too

Only problem: this is super slow

# Much better idea

- Building an image is basically just running containers
- Use OverlayFS to make the image writable
- This lets us build content indexes very easily
- "Only" need to implement a Dockerfile parser

```rust
Instruction::Run(ins) => {
    use dockerfile_parser::ShellOrExecExpr::*;
    let mut run_env = env.clone();
    run_env.extend(self.task_env.clone()); // User environment overrides image variables.
    match &ins.expr {
        Shell(cmd) => {
            let cmd = cmd.to_string();
            let mut args: Vec<&str> =
                self.shell.iter().map(|s| s.as_ref()).collect::<Vec<_>>();
            args.push(&cmd);
            runc.exec(&args, &format!("{cmd:?}"), work_dir, &run_env)?
        }
        Exec(args) => {
            let args = args.elements.iter().map(|s| s.as_ref()).collect::<Vec<_>>();
            runc.exec(&args, &format!("{args:?}"), work_dir, &run_env)?
        }
    }
}
```

# What about scheduling?
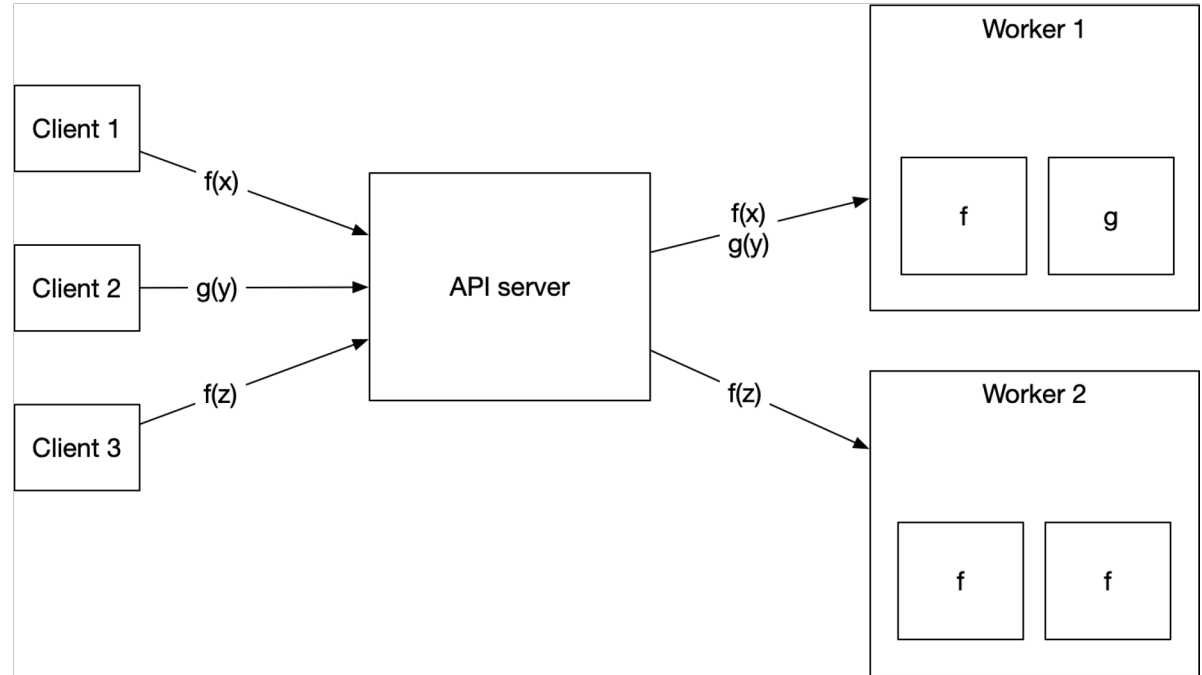
What did we build so far:

- ~~Run custom images very fast~~
- ~~Build custom images very fast~~
- Maintain a pool of worker instances
- Allocate jobs to workers
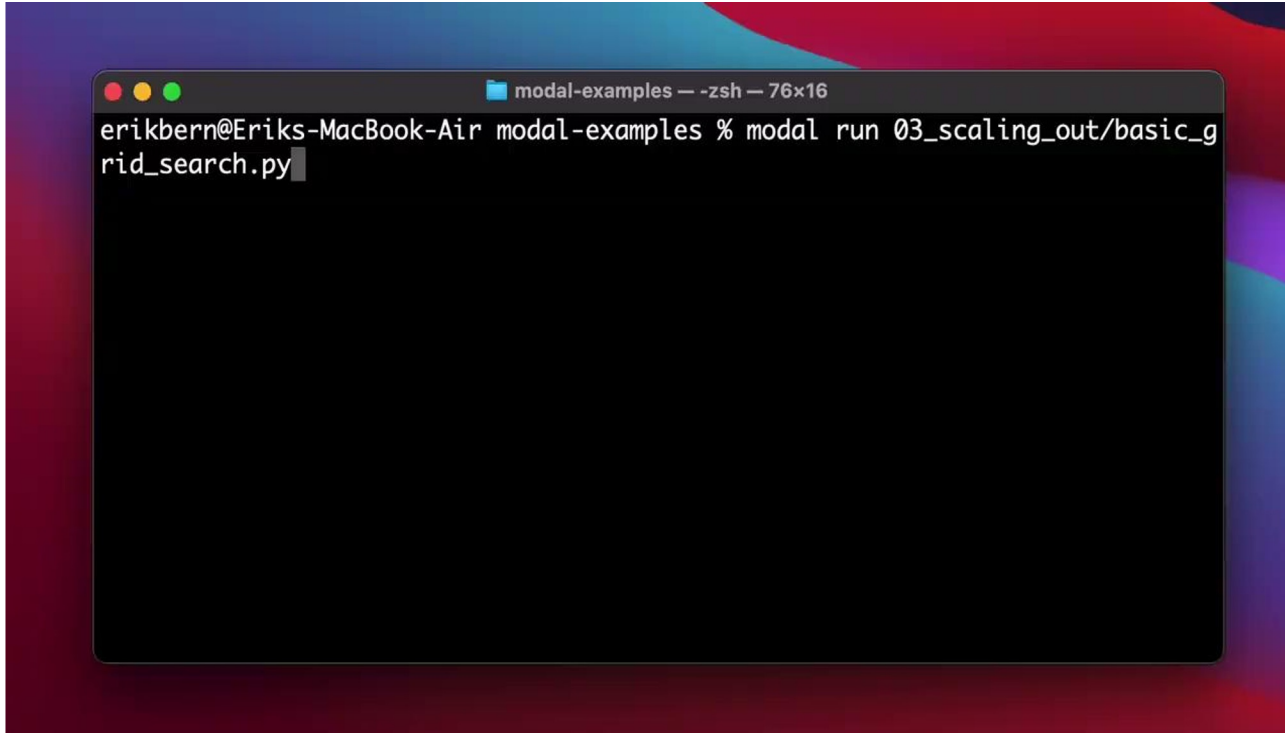
# Let's run our own resource pool

- Launch & terminate instances on AWS & GCP
- We can launch an instance in about 40s
- "Overprovision" so we always have a bit of spare capacity
- We benefit from multi-tenancy
- Every worker reports available CPU & memory every 2s

# Turning this into a function-as-a-service platform

- Main trick: reuse the same container for multiple function calls
- Autoscale on-demand, scale down to zero quickly
- Super useful for GPUs
- Need fast cold start

# What does this let us do?

# What are some use cases?

- Lots of Stable Diffusion and Dreambooth
- Also computational biotech, web scraping, data pipelines, and many other things



Photo you provide

Headshot examples you get back

# Was it dumb to build this in-house?

Maybe? But

- Docker is too slow & limited for what we needed
- It would have taken too much work getting Kubernetes to do this
- AWS Lambda is too expensive and limited

# Thanks!

Questions?
erik@modal.com
@bernhardsson