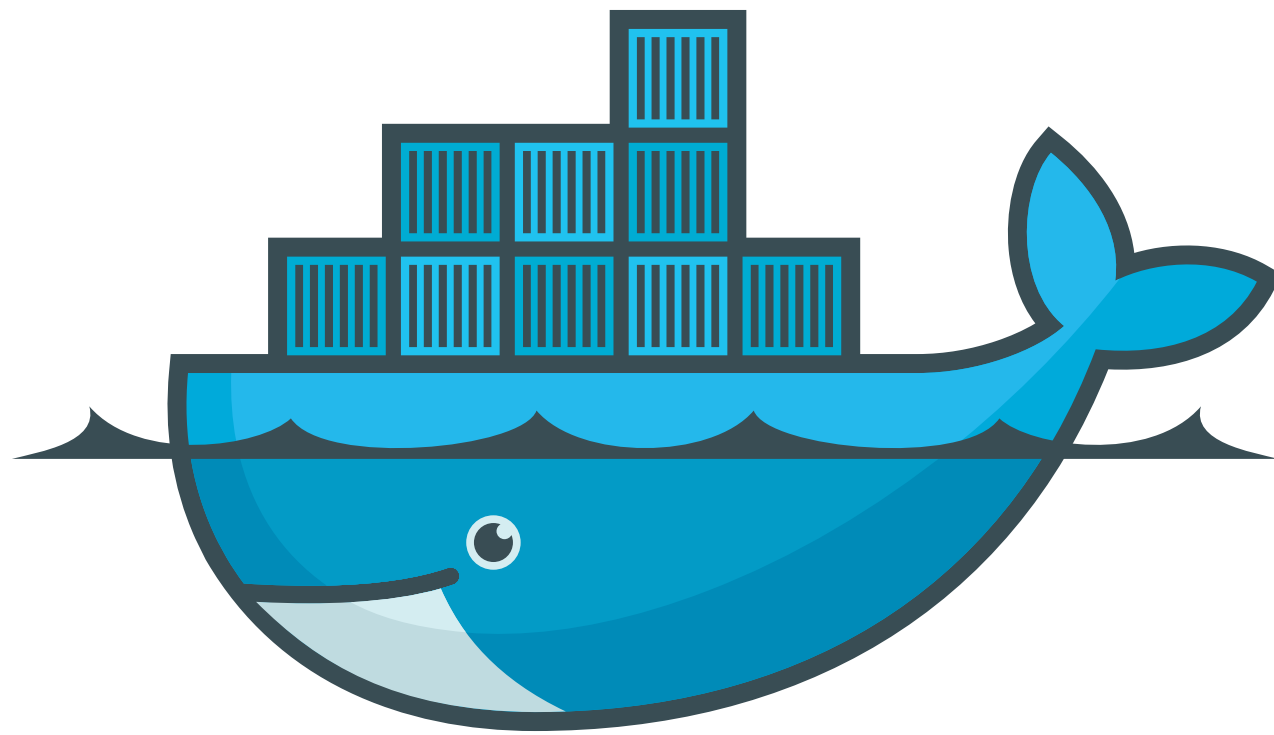# Docker-encapsulated pipelines

## Brief overview of my summer placement

Stefan Dang · Wellcome Trust Sanger Institute · Hinxton, 14.10.14

# Introduction
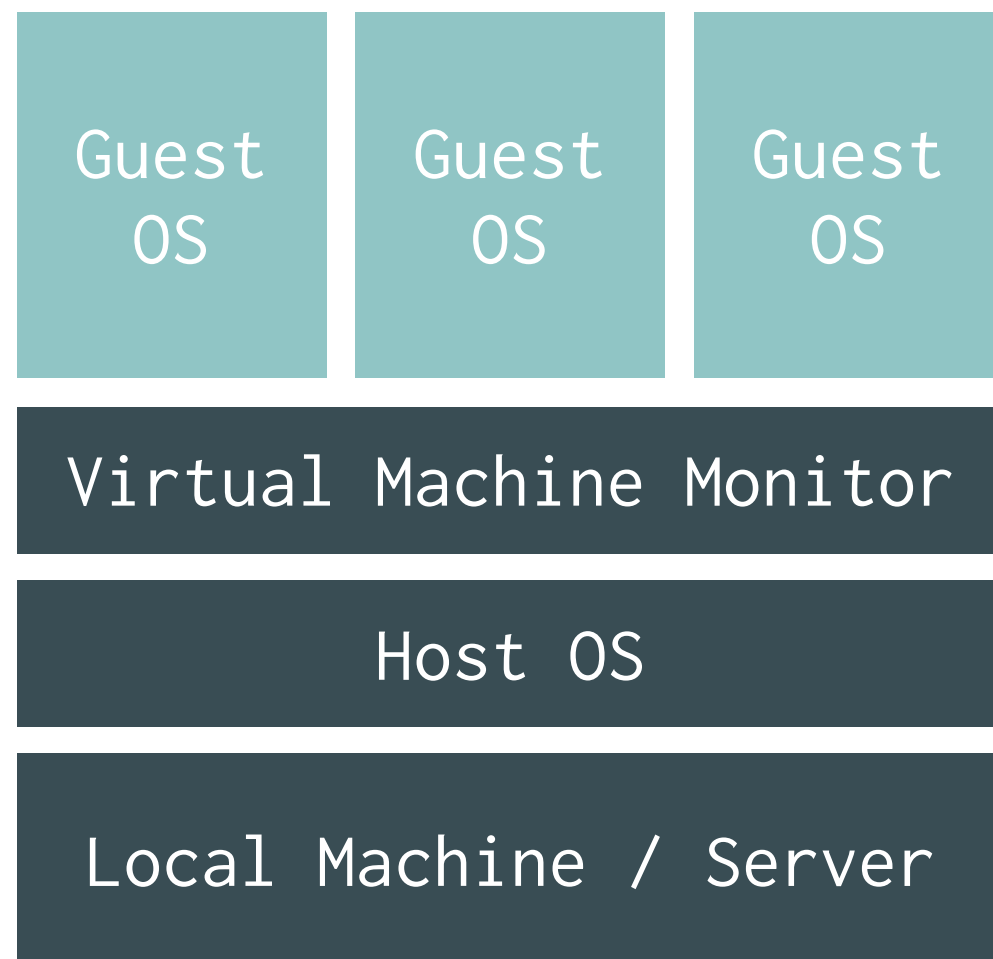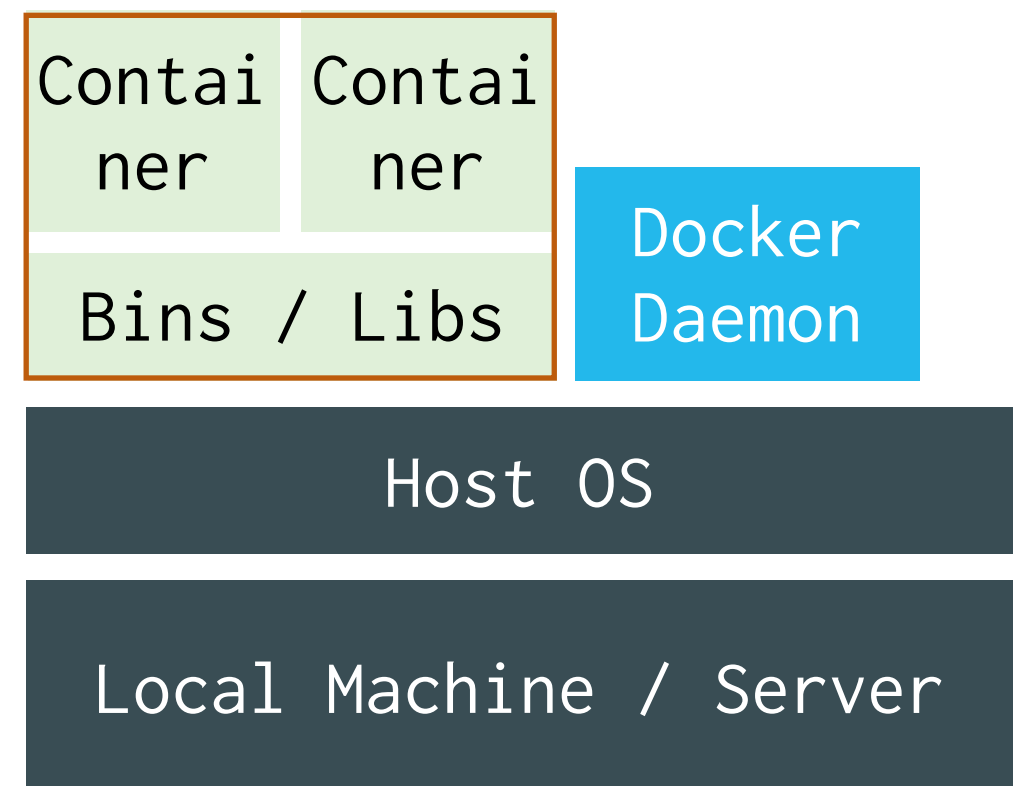
docker

docker.com

# Virtualization

## Virtual Machine

| Guest OS | Guest OS | Guest OS |
|---|---|---|

| Virtual Machine Monitor |
|---|

| Host OS |
|---|

| Local Machine / Server |
|---|

## Linux Containers

| Contai ner | Contai ner |
|---|---|
| Bins / Libs | |

| Docker Daemon |
|---|

| Host OS |
|---|

| Local Machine / Server |
|---|

# Docker in Detail

## Underlying Infrastructure

| libcontainer / LXC |
| LayerFS | cgroup | Namespaces |
| Linux kernel |
| Local Machine / Server |

## Client-Server-Application

Client

Daemon

Host

Registry

# Dockerfile

```
FROM debian:jessie
MAINTAINER Stefan Dang <sd15@sanger.ac.uk>

ENV FOO_VERSION 0.6.7

# Download dependencies & sources
RUN apt-get update && sudo apt-get install -yqq [Dependencies]

# Build foo
RUN git clone -b $FOO_VERSION https://github.com/foo.git
WORKDIR ./foo
RUN ./configure && make && make install

COPY ./entrypoint.sh /entrypoint.sh
ENTRYPOINT ["/entrypoint.sh"]
```
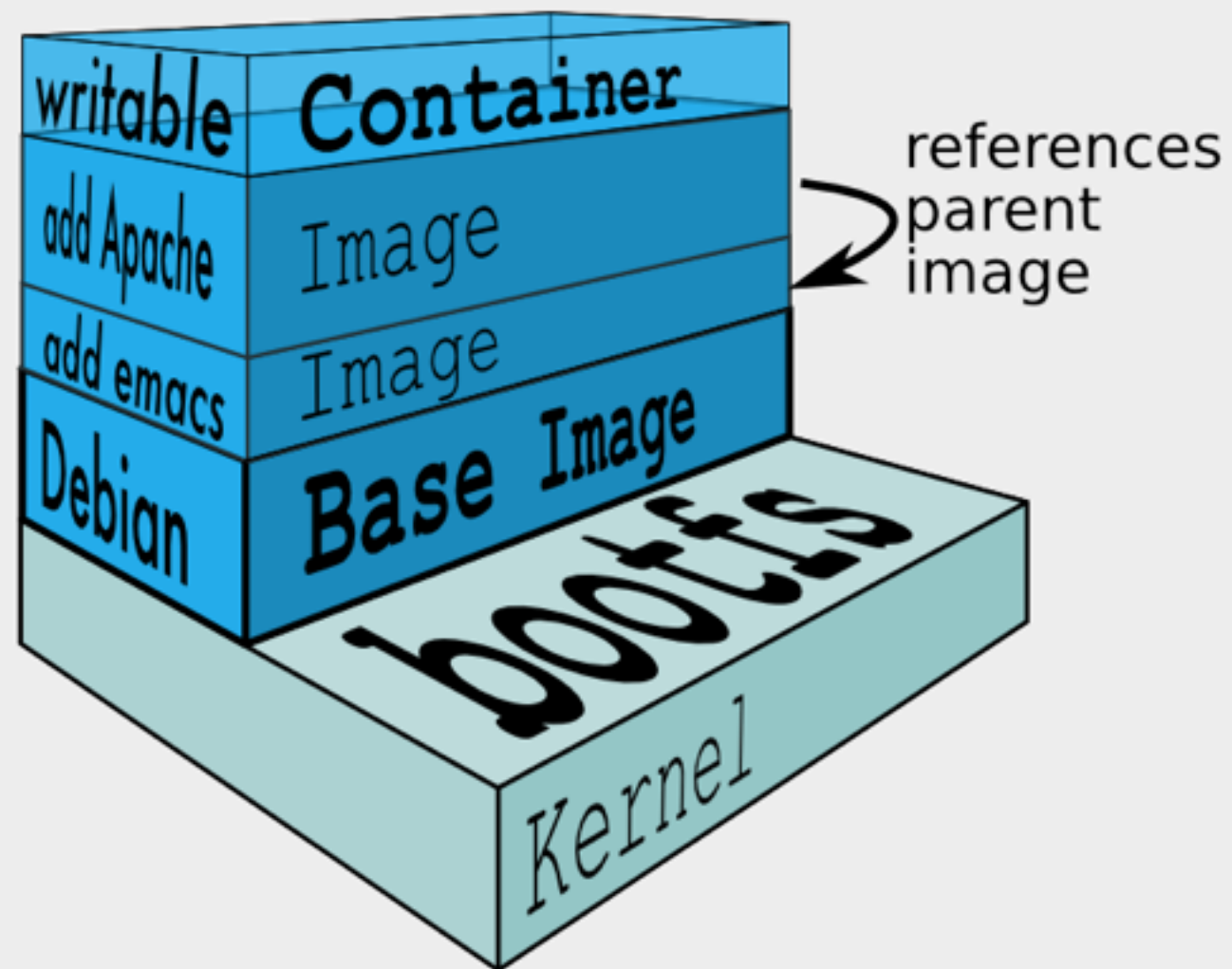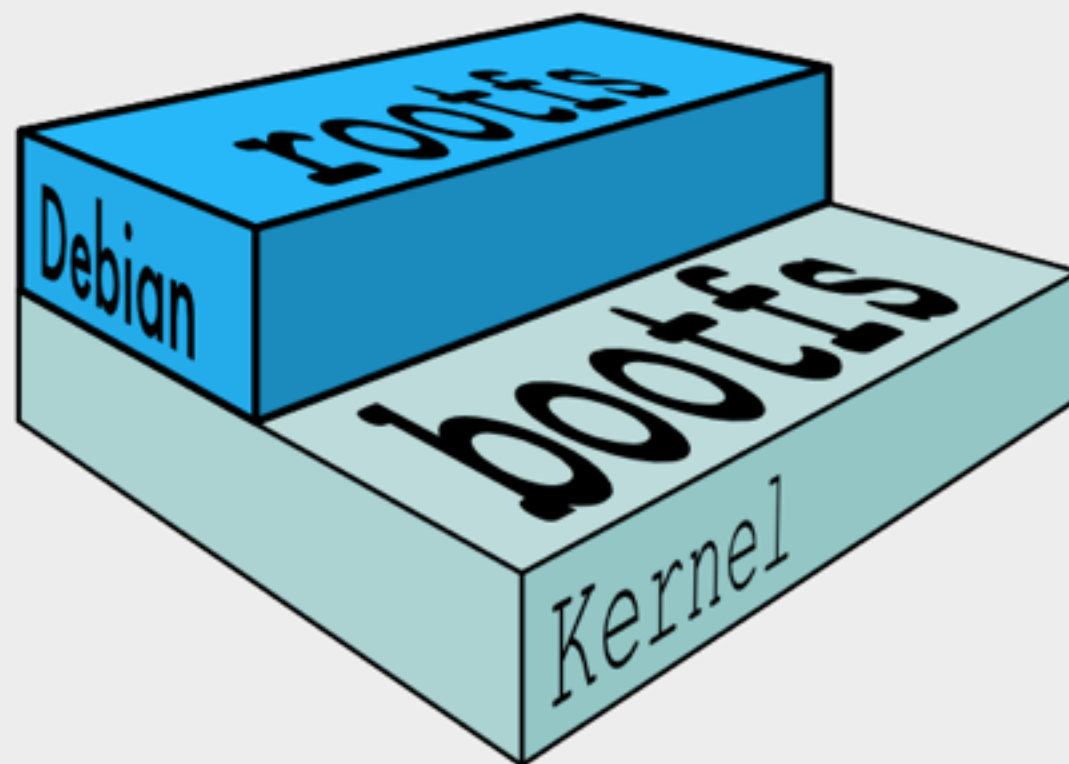
```
$ docker build -t ImageName ./DockefileDir/
```
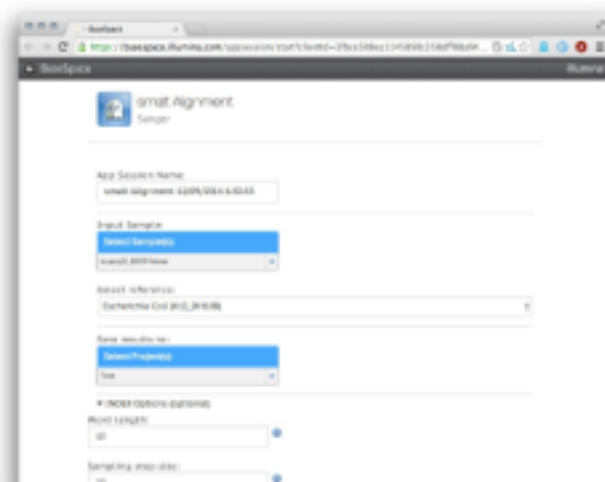
references
parent
image

```
$ docker run -it ImageName
```

# Use Cases

# Illumina BaseSpace Native Apps



**Input Form**

JSON

Javascript

**App / Pipeline**

AWS EC2

**Output Form**

dotfluid

# Smalt Docker Container

Input Form

entrypoint.sh

Flow: smalt + Post-Processing

iGenomes

Input

Output

Output Report

# Smalt Flow (gist)

```
$ smalt index

$ smalt map \
  > bamsort \
  | bamstreamingmarkduplicates \
  | tee >(samtools flagstat) >(samtools stats) \
  | bamrecompress md5=1 index=1 \
  > "out.bam"

$ plot-bamstats
```

**Dockerfile**

```dockerfile
1   FROM debian:jessie
2   MAINTAINER Stefan Dang <sd15@sanger.ac.uk>
3
4   # Bioinformatics Tools Versions
5   ENV bambamc_version 0.0.50-release-20140430085950
6   ENV smalt_version 0.7.6
7   ENV samtools_version 1.1
8   ENV libmaus_version libmaus_experimental_0_0_153
9   ENV biobambam_version biobambam_experimental_0_0_163
10  # TODO: smalt and samtools still have hard-coded WORKDIR instructions, see
11
12  # Install dependencies
13  RUN apt-get update -q &&\
14      apt-get install -qy build-essential \
15                          autoconf \
16                          automake \
17                          git \
18                          gnuplot \
19                          libtool \
20                          libncurses5-dev \
21                          libncursesw5-dev \
22                          pkg-config \
23                          wget \
24                          zlib1g-dev
25
26
27  # Build bioinformatics tools in /home:
28  WORKDIR /home
29
30  # Build bambamc lib (for smalt bam support)
31  RUN git clone -b $bambamc_version https://github.com/gt1/bambamc.git bamba
32  WORKDIR ./bambamc
33  RUN autoreconf -i -f &&\
34      ./configure &&\
35      make && make install
36  WORKDIR ..|
37
38  # Build smalt
39  RUN wget -qO- http://sourceforge.net/projects/smalt/files/smalt-$smalt_ver
40      | tar -xz
41  # TODO: Use env variable as soon as docker 1.3 is released
42  WORKDIR ./smalt-0.7.6
```

**smalt_entrypoint.sh**

```bash
11  set -o pipefail
12  set -e
13
14  # Globals
15  INDEX=$1; shift                    # $1
16  PROJECT_ID=$1; shift               # $2
17  INDEX_WORDLEN=$1; shift            # $3
18  INDEX_STEPSIZE=$1; shift           # $4
19  INSERT_MAX=$1; shift               # $5
20  INSERT_MIN=$1; shift               # $6
21  COUNTER=0                          # Make sure alignment has run at least o
22
23  # Catch empty input, set to standard values
24  [[ -z "$INDEX_WORDLEN" ]] && INDEX_WORDLEN="13"
25  [[ -z "$INDEX_STEPSIZE" ]] && INDEX_STEPSIZE="$INDEX_WORDLEN"
26  [[ -z "$INSERT_MAX" ]] && INSERT_MAX="500"
27  [[ -z "$INSERT_MIN" ]] &&  INSERT_MIN="0"
28
29  # Indexing
30  smalt index -k "$INDEX_WORDLEN" -s "$INDEX_STEPSIZE" "$INDEX" "$INDEX.fa"
31
32  # Prepare output folders, respecting Basespace naming convention
33  mkdir -p "/data/output/appresults/$PROJECT_ID/smalt"
34
35  # Iterate over all files
36  for input_file in /data/input/samples/*/*; do
37    filename=$(basename "$input_file" .fastq.gz)
38    output_file=/data/output/appresults/$PROJECT_ID/smalt/$filename
39
40    # Only process R1 (following Illumina naming convention), check for R2 b
41    if [[ $filename =~ _R1_[0-9]{3} ]]; then
42      gzip -dc "$input_file" > input.fastq || err "Could not decompress $fil
43
44      # Set post-processing pipeline:
45      # bamsort | bamstreamingduplicates | samtools flagstat & stats | recom
46      mkfifo postproc_pipe && \
47      bamsort level=0 SO=coordinates fixmates=1 adddupmarksupport=1 \
48      < postproc_pipe \
49      | bamstreamingmarkduplicates level=0 \
50      | tee >(samtools flagstat - > "$output_file.flagstat") \
51           >(samtools stats - > "$output_file.stats") \
52      | bamrecompress md5=1 md5filename="$output_file.md5" \
```
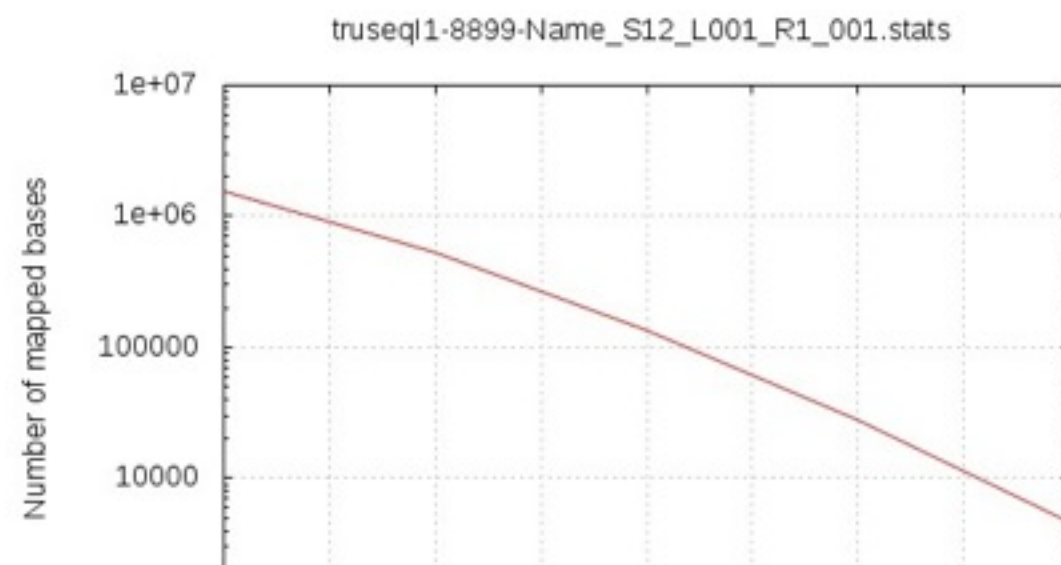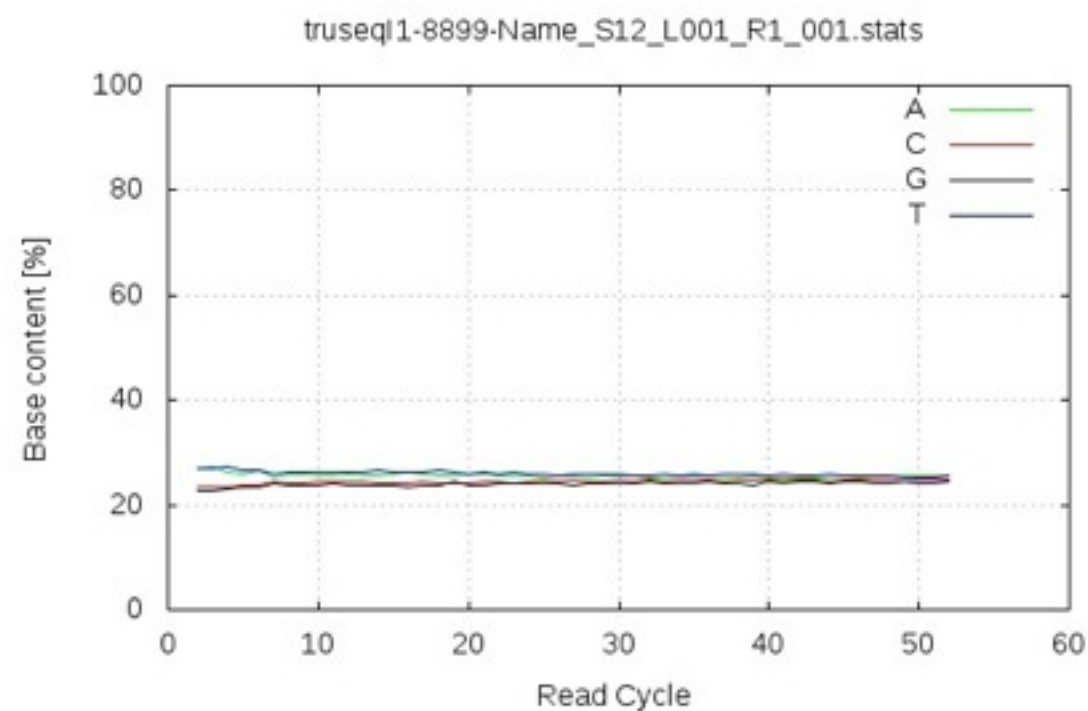
⬇ Download Analysis  ✎ Rename analysis  Move to Trash

View Trash 🗑

ⓘ Analysis Info

⇥ Inputs

▤ Output Files
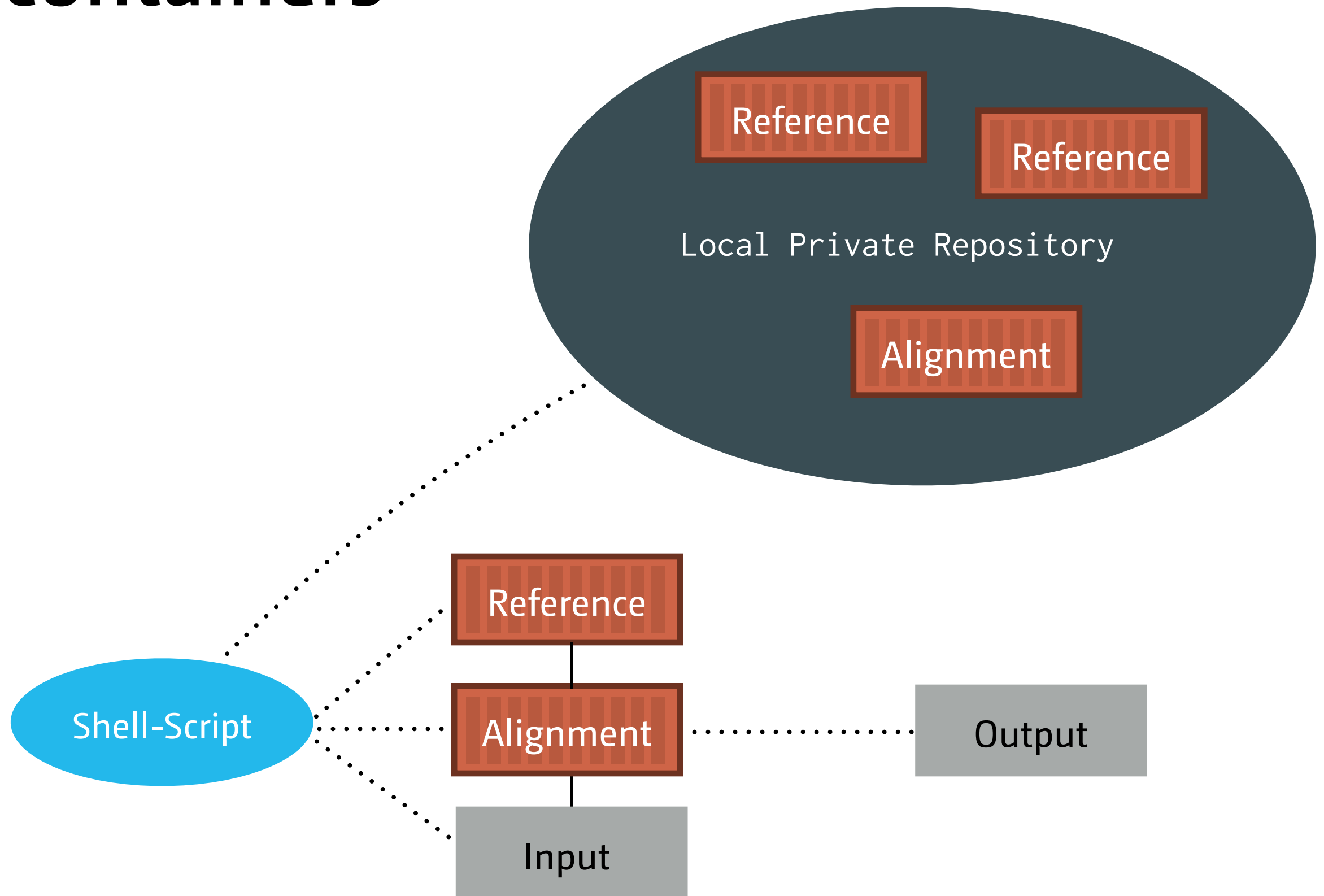
📊 Analysis Reports

smalt

## FLAGSTATS

65844 + 0 in total (QC-passed reads + QC-failed reads) 838 + 0 duplicates 61256 + 0 mapped (93.03%:-nan%) 65844 + 0 paired in sequencing 32922 + 0 read1 3292
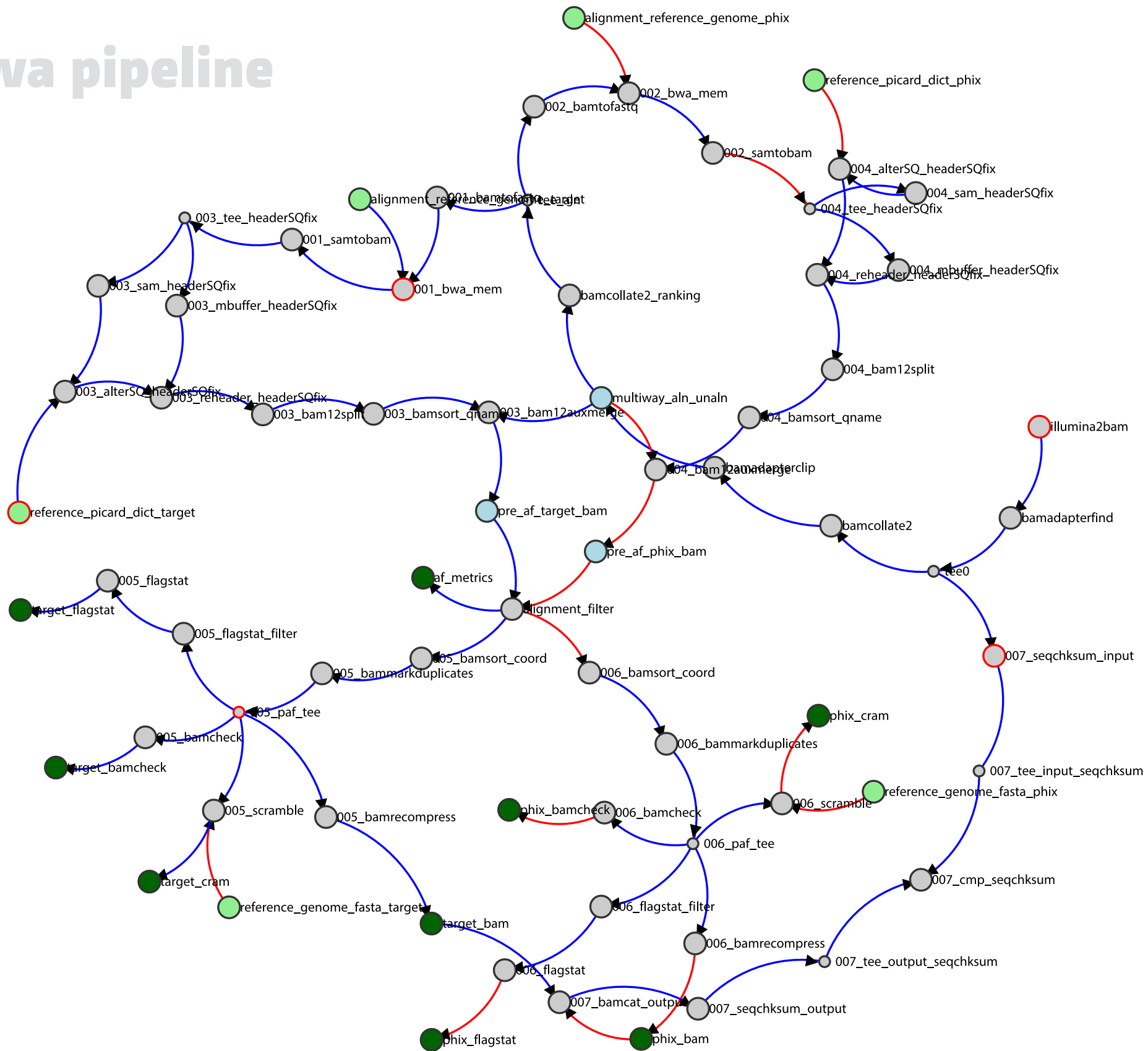
## PLOT FLAGSTATS

truseql1-8899-Name_S12_L001_R1_001.stats



truseql1-8899-Name_S12_L001_R1_001.stats

# Pipelines in orchestrated Docker containers

# p4 bwa pipeline

libmaus
biobambam

bwa

illumina2
bam

picard

samtools

# Docker Build System

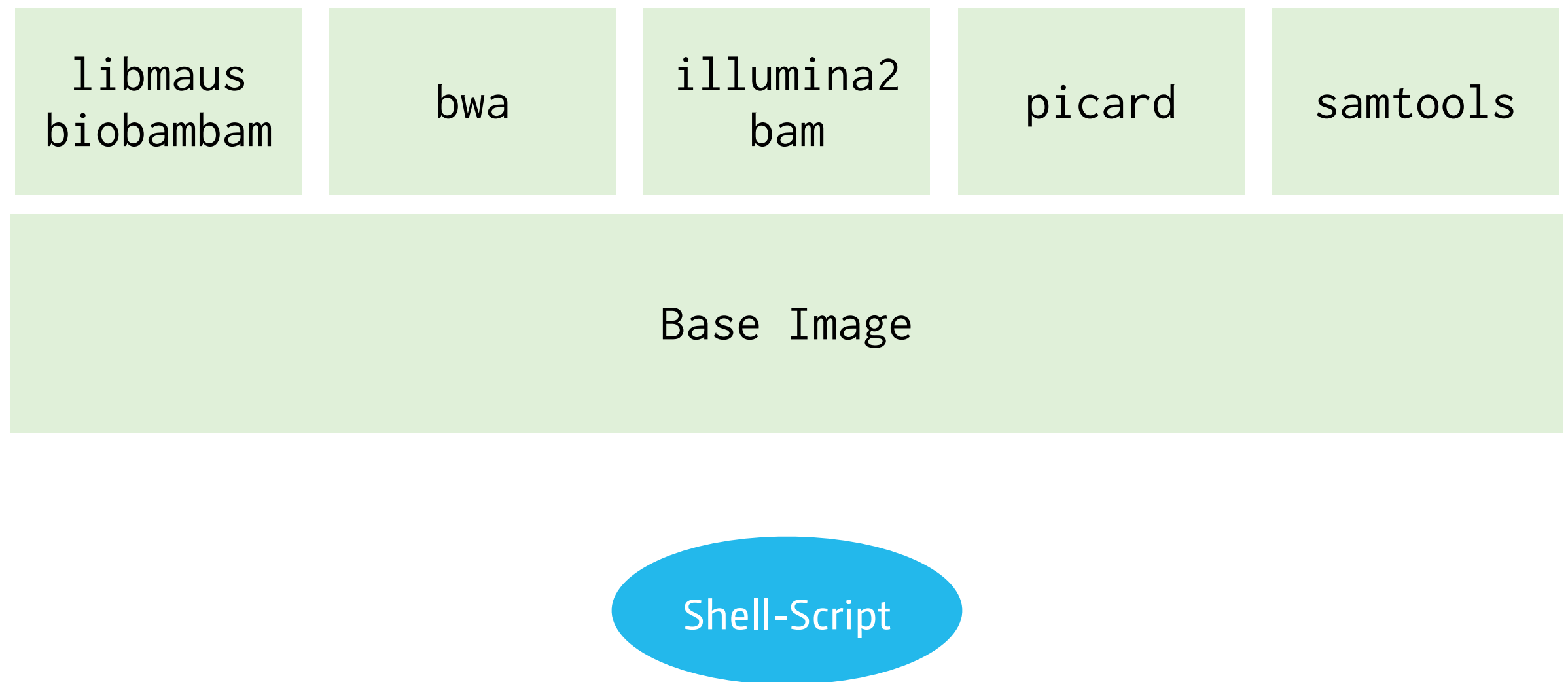| | | | | |
|---|---|---|---|---|
| libmaus biobambam | bwa | illumina2 bam | picard | samtools |

Base Image

Shell-Script

# Overhead

```
$ time docker run ubuntu:14.04
0.01s user 0.02s system 2% cpu 1.007 total

$ docker images

REPOSITORY                          VIRTUAL SIZE
ubuntu:14.04                         194.9 MB
debian:jessie                        120.0 MB
p4                                  1259.0 MB
p4_flattened                         484.1 MB
p4_build-system                      294.2 MB
PhiX_iGenome                          16.9 MB
```

# Summary

# Drawbacks

## System requirements

64bit, Linux 3.8 or later kernel, Storage Driver Support

Boot2Docker

## Security

docker daemon requires root privileges

less mature / tested compared to VMs

less isolation compared to VMs

# Advantages

## Open Source

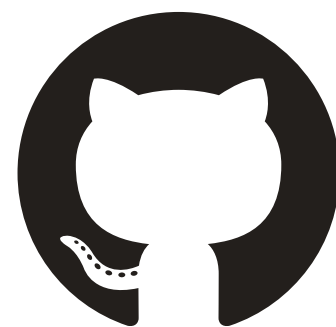github, Go, Apache 2.0 license

## Speed & Maintenance

predictable

repeatable

managed

# Thanks

Marina Gourtovaia

Kevin Lewis

David Jackson

# Thank you!

# Questions?