

# 9. The Cloud & Data

CSCI 2541 Database Systems & Team Projects

Gabe modified from Wood

What is the oldest piece  
of software you  
remember using?

# Software Changed



**Then**



**Now**

Where and how we run programs has changed

- Network connected
- Mobile
- Multi-media content
- Shared by lots of users

# Cloudy Buzz

*Mobile*

!

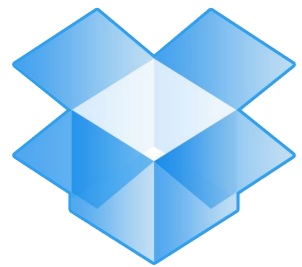


Google Docs



iCloud

To the cloud!



Dropbox

flickr™

*Fast!*

XBOX  
LIVE



amazon  
web services™

*Free\*!*

*Powerful!*



# What *is* a cloud?

<spoiler alert>

It's not in the sky

it's not made of water droplets

</spoiler alert>

# Some big buildings...



Microsoft's Dublin data center

# ...and computers...

## Giant warehouses

- The size of 10 football fields
- 10s of thousands of servers
- Petabytes of storage



# ...interconnected...

**Level(3)**<sup>TM</sup>  
COMMUNICATIONS



# ...around the world...



## Undersea Cables

- Connect all continents except Antarctica
- First deployed in 1850s



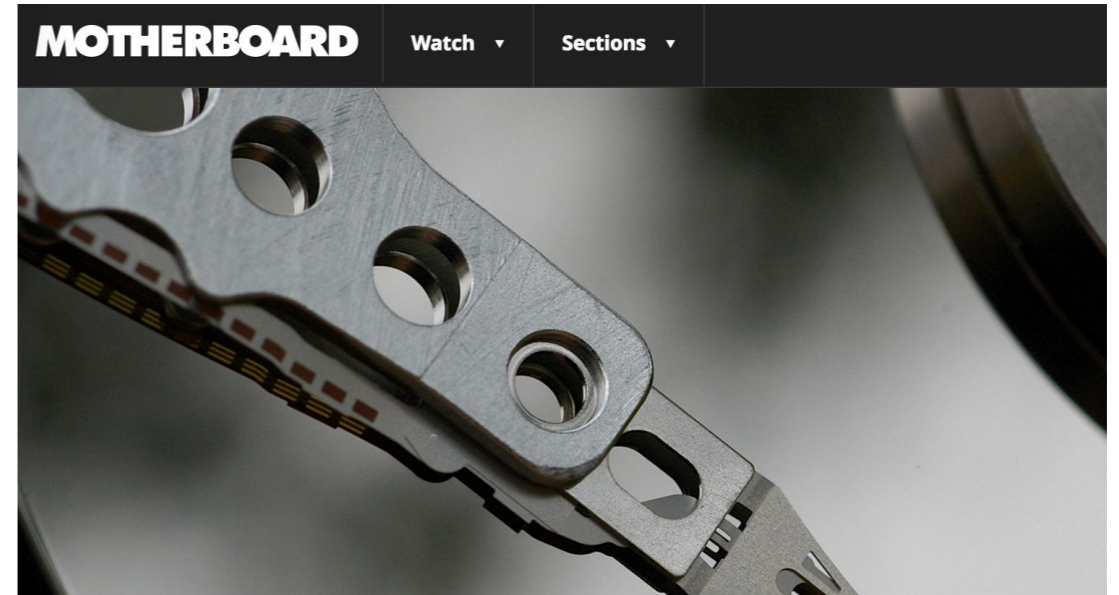
<http://www.cyprusupdates.com>

m

# ...that break a lot.



Lightning causes Amazon outages (2009 and 2011)



**A Loud Sound Just Shut Down a Bank's Data Center for 10 Hours**

September 11, 2016 // 02:00 PM EST



Anchor hits underwater Internet cable (Feb 2012)



Comcast down after hunter shoots cable (2008)

Or if you're really  
unlucky...



**VS**



# *Cloud* Defined

**cloud:** */klaʊd/* noun

A **large** collection of computers, accessible over a **network**, running many different types of software as a **shared** service

Must be:

efficient, scalable, secure, reliable, *elastic*

# Cloud Examples

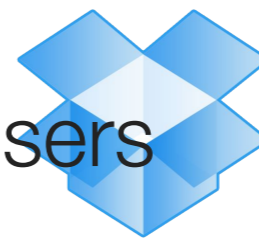


Shared, worldwide infrastructure to host email services for many users and organizations

- ~900,000 servers in 2014

Shared storage service

- ~10,000 servers and 200 million users in 2013



**Dropbox**

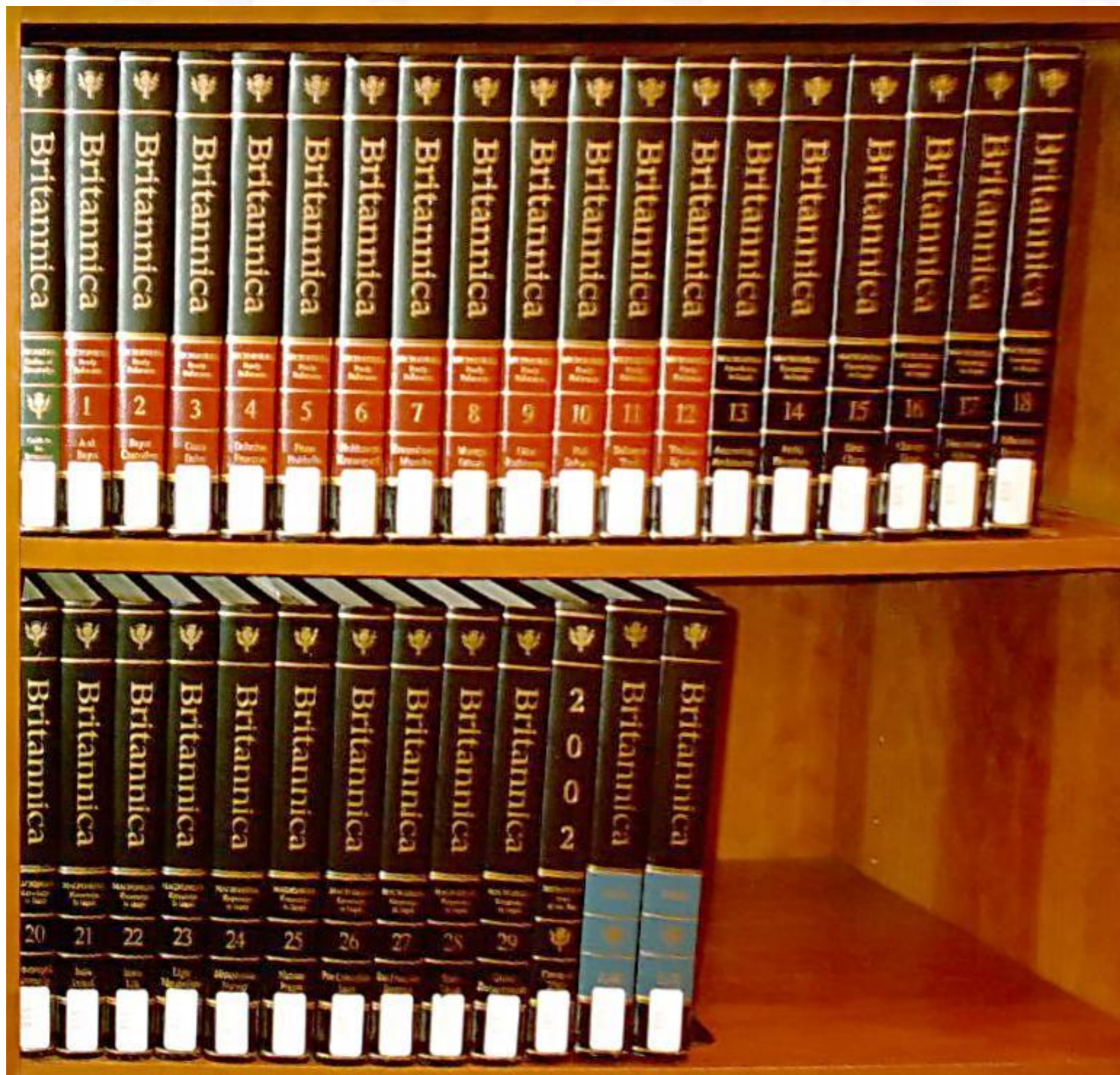


Shared computing infrastructure that developers, companies, and students can easily get access to

- ~1.4 million servers in 2014

Why do we need all of  
this physical  
infrastructure?

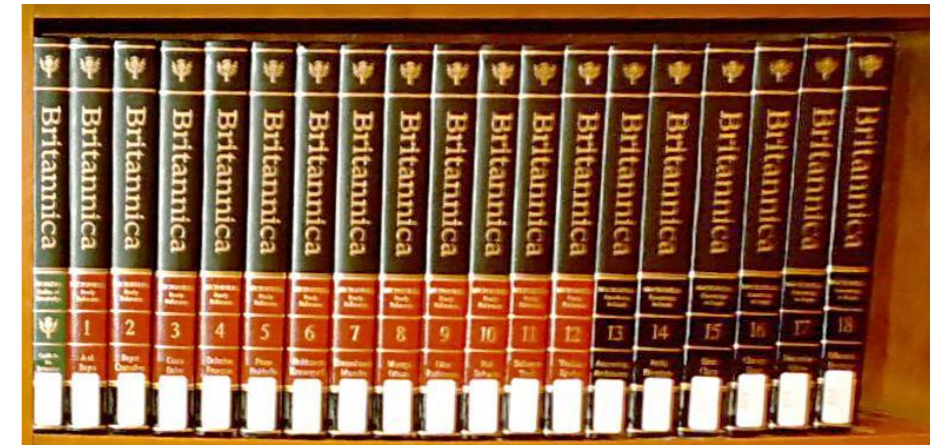
# What is this???



# Encyclopedias

## Encyclopedia Britannica

- 40,000+ articles
- 32 hard bound volumes (32,640 pages)



## Microsoft Encarta

- 60,000+ articles
- 1 CD-ROM (**700 MB**)



## Wikipedia

- 6,383,000 articles (in English)
- More than **5 TB** of text (about 7,500 CDs)



# Mega whats?

700MB vs 5TB

Mega	Million	$1024 \times 1024 =$ $\sim 1,000,000$
Giga	Billion	$1024 \times 1024 \times 1024 =$ $\sim 1,000,000,000$
<b>Tera</b>	<b>Trillion</b>	$1024 \times 1024 \times 1024 \times 1024 =$ <b><math>\sim 1,000,000,000,000</math></b>

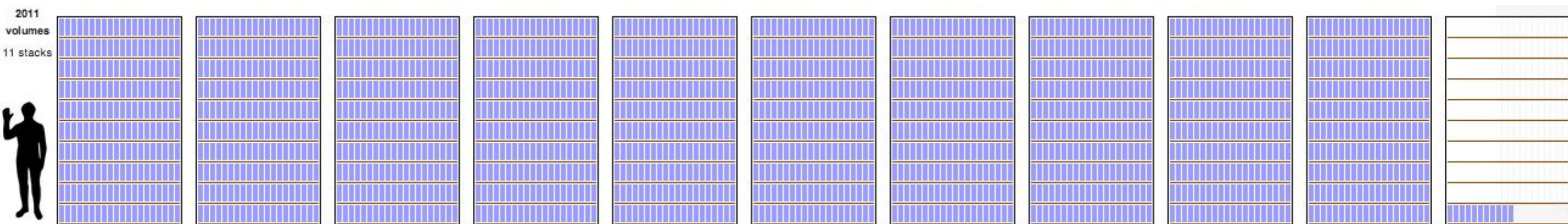
200 photos vs 1.4 million photos

# Encyclopedias

Wikipedia... in print

- ~~1,763 volumes~~
- (no, this does not exist)

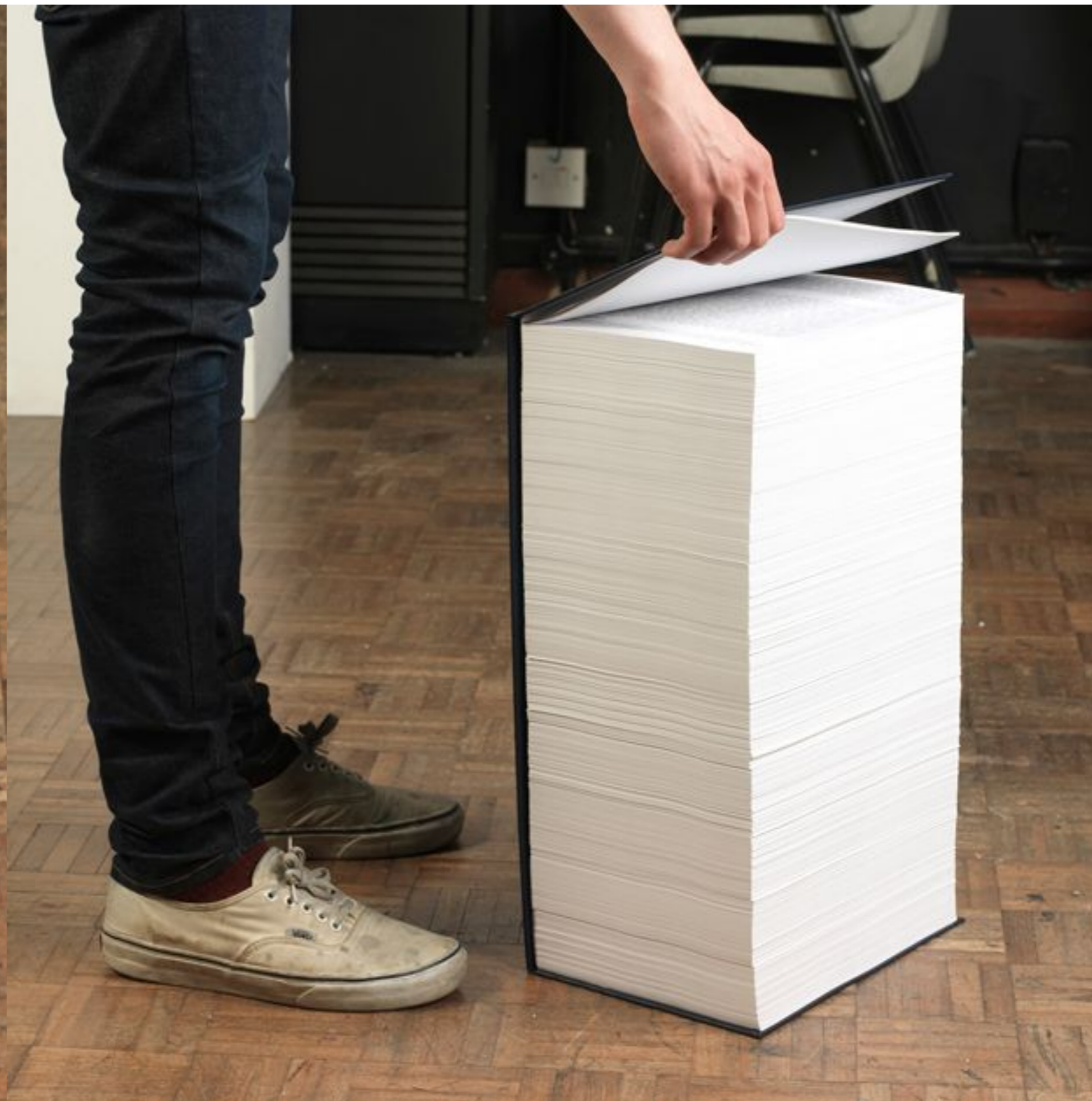
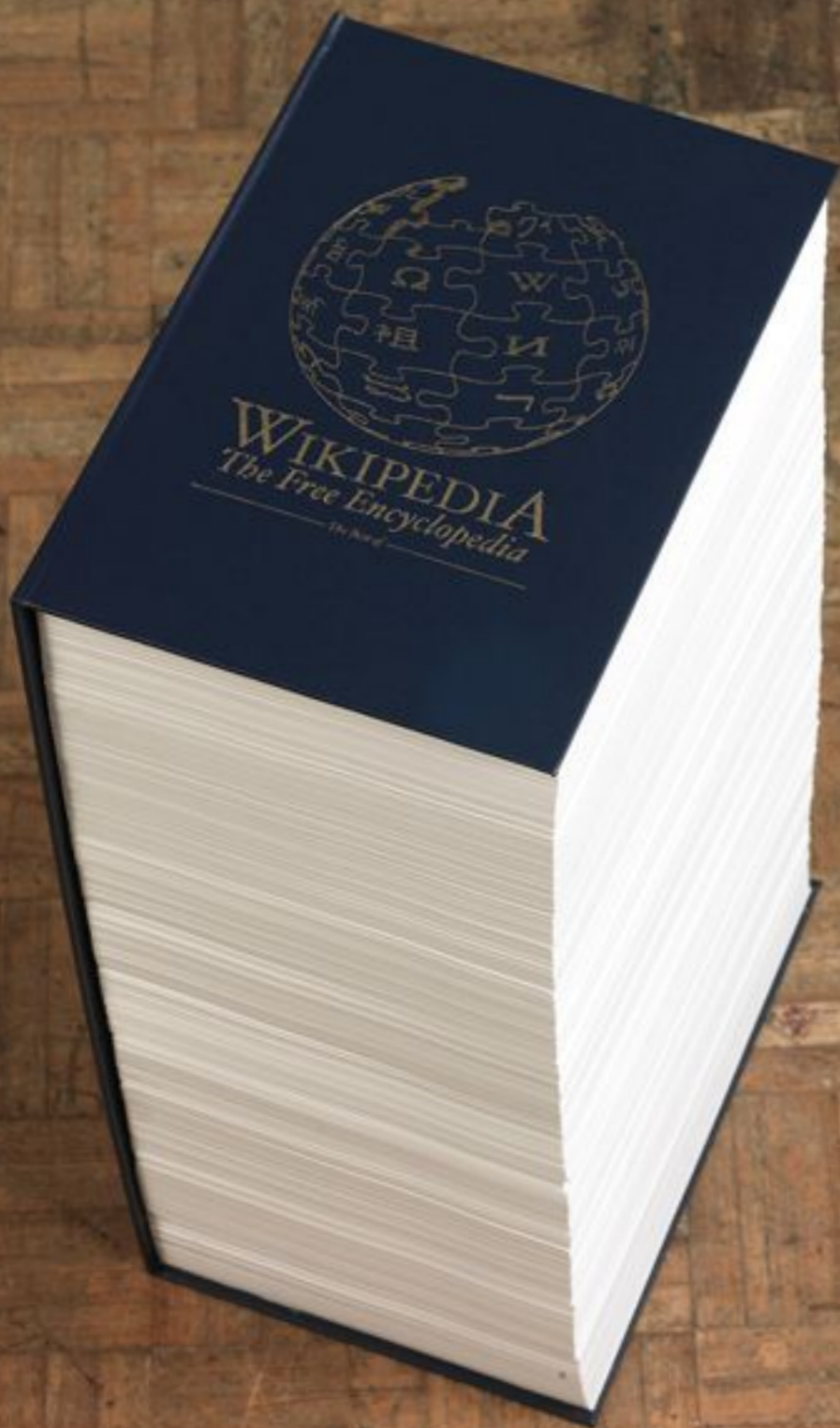
Now grown to **3,024** volumes  
and >30TB of data!



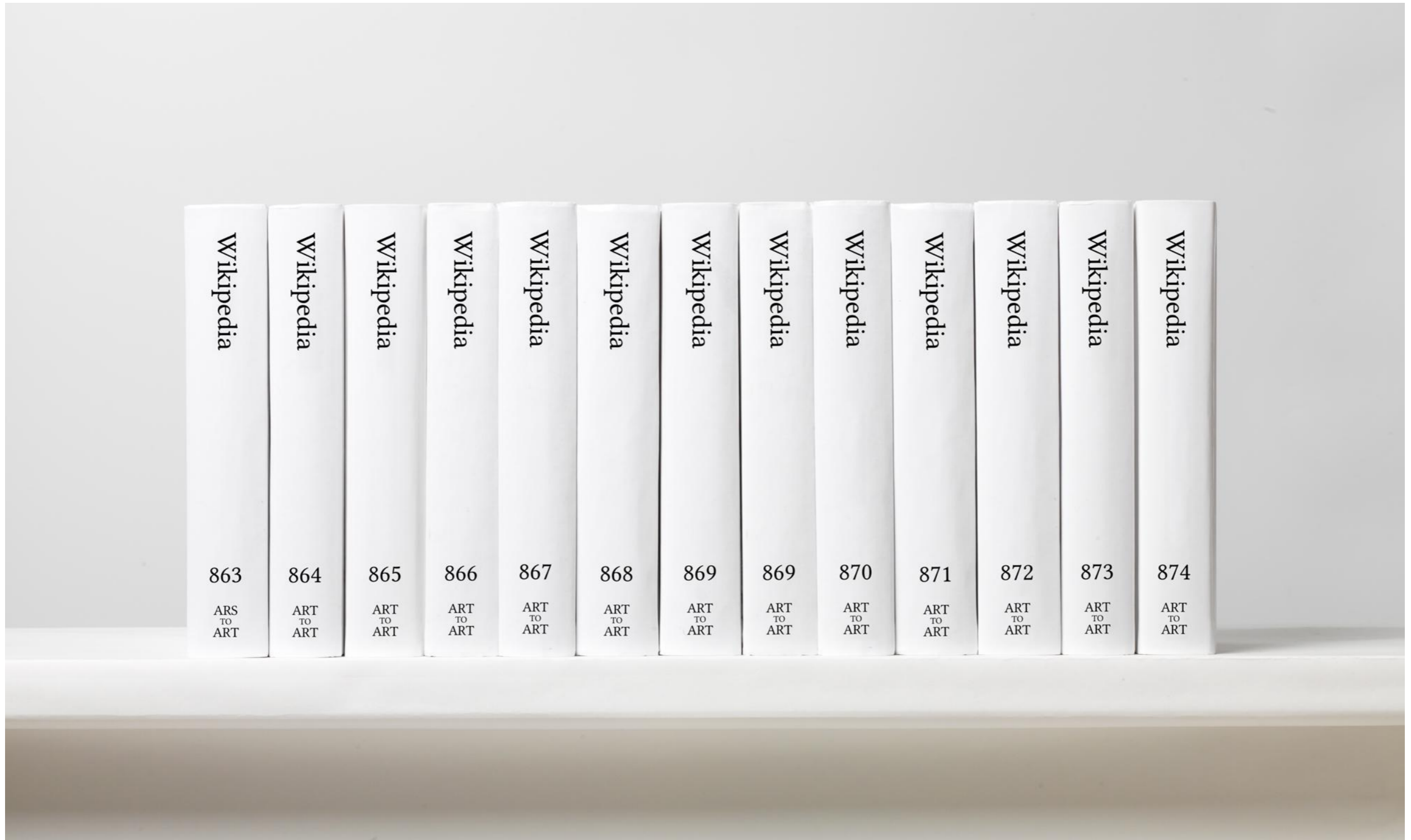
[http://en.wikipedia.org/wiki/Wikipedia:Size\\_in\\_volume](http://en.wikipedia.org/wiki/Wikipedia:Size_in_volume)

S

# 0.01% of Wikipedia



# It exists! (sort of)



# Own it!

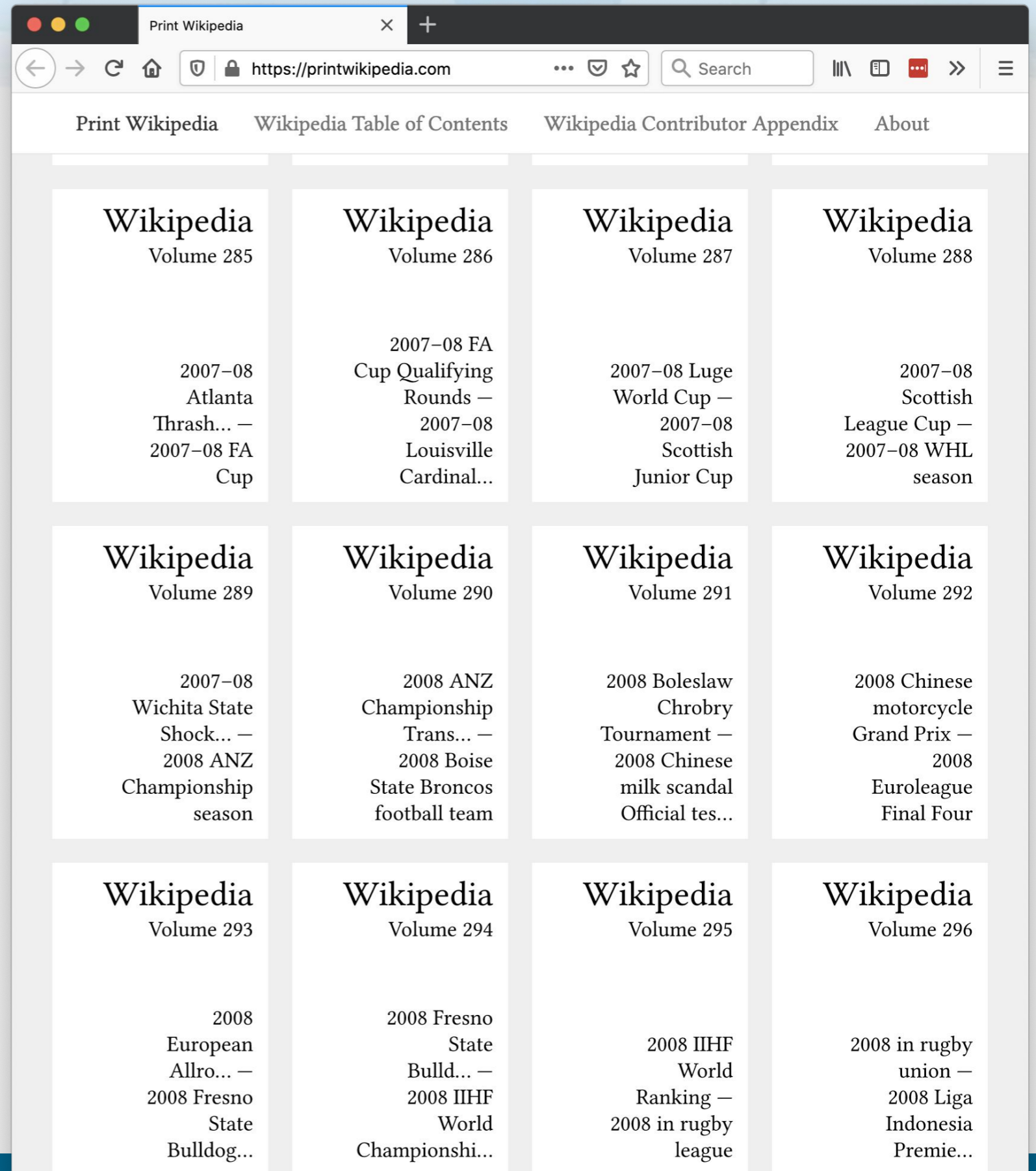
## Just \$80\*!!!

\*per volume

## 7,473 volumes each with 700 pages

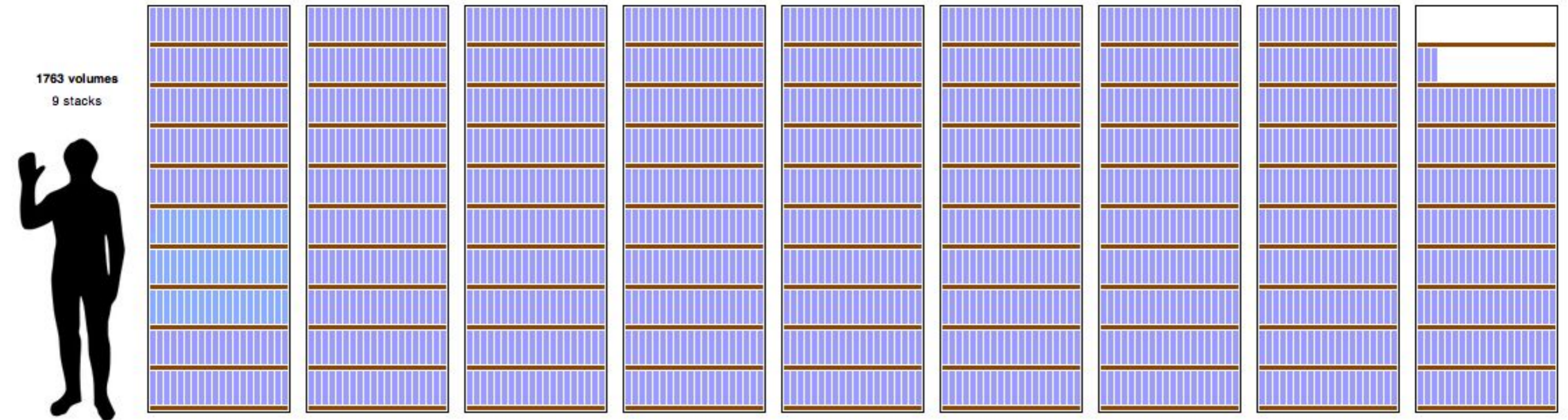
## Print on demand

## <https://printwikipedia.com>



# Big Data in Perspective

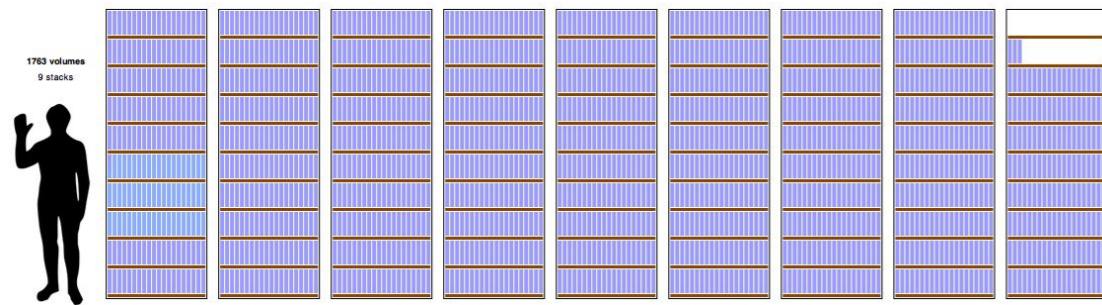
**Wikipedia** - 5TB of text



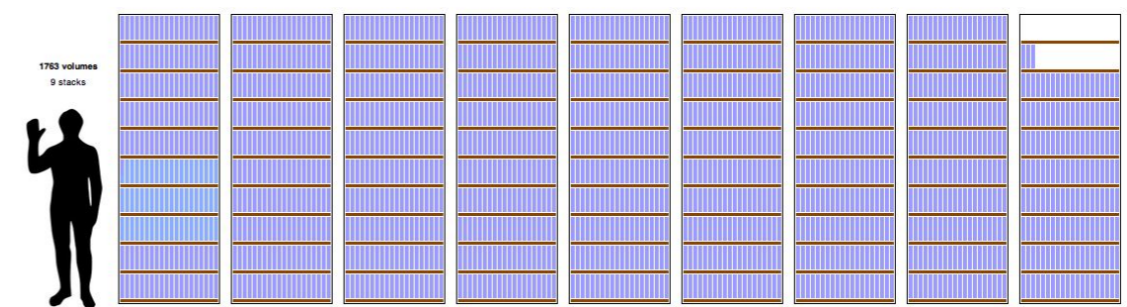
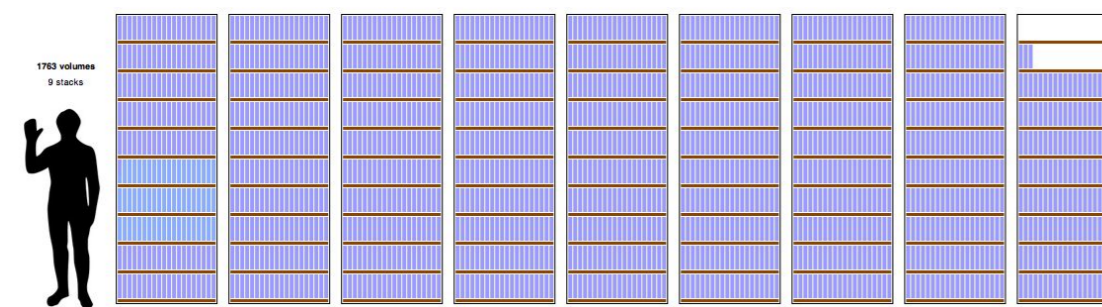
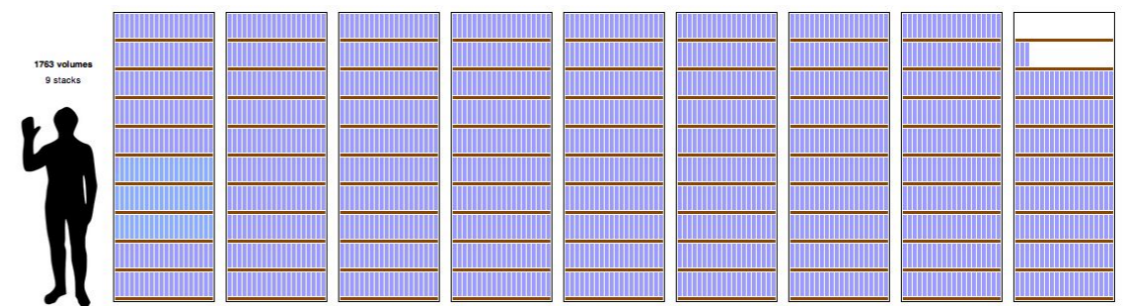
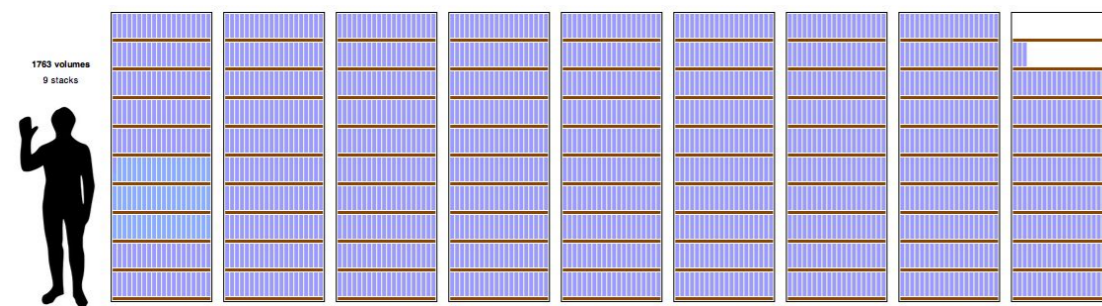
**Facebook** - ???

# Big Data in Perspective

**Wikipedia** - 5TB of text



**Facebook** - 20TB of photos added *each week*



# Big Data in Perspective

A large grid of small human icons representing data volume. The grid is composed of many small human silhouettes arranged in a regular pattern, filling most of the slide area. In the top-left corner, there is a small icon of a person standing next to a grid of colored squares (purple, blue, green, yellow, red) representing data storage or processing units.

**Facebook** - 1,000TB of photos added *per year*

# Big Data in Perspective

**Facebook** - 1,000TB of photos added *per year*

**Google** - 20,000TB of data processed *per day*

# Big Data in Perspective

**Google** - 20,000TB of data  
processed *per day* - *in 2008*

# Big Data in Perspective

**Google** - 20,000TB of data processed *per day* - *in 2008*

**Google** - Estimated 200,000TB of data processed *per day* - *in 2018*

40,000  
wikipedias per  
day!

How can google  
process *so much*  
information *so*  
*quickly?*

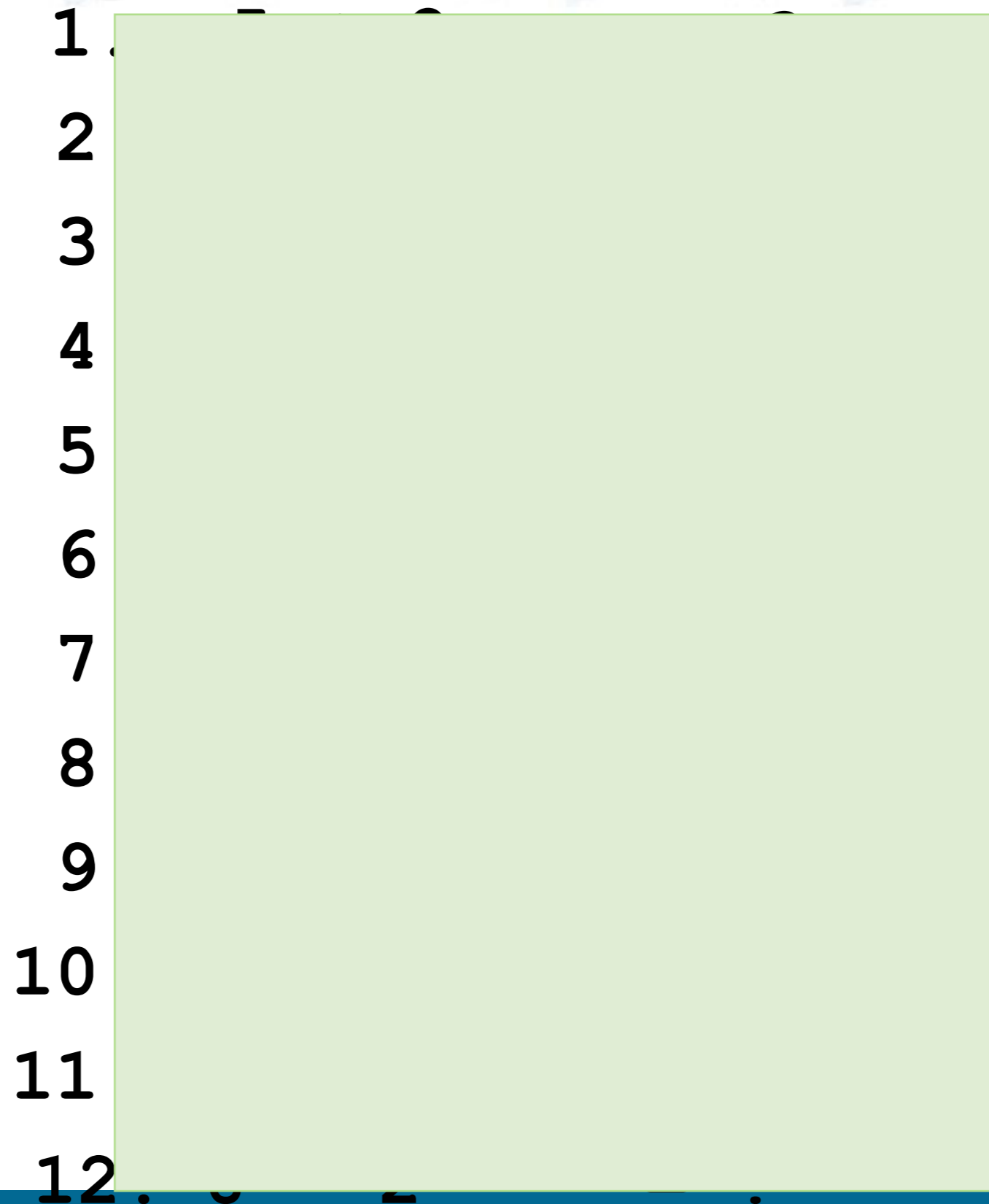
# Processing Data Quickly

1.  $3 + 6 = ?$

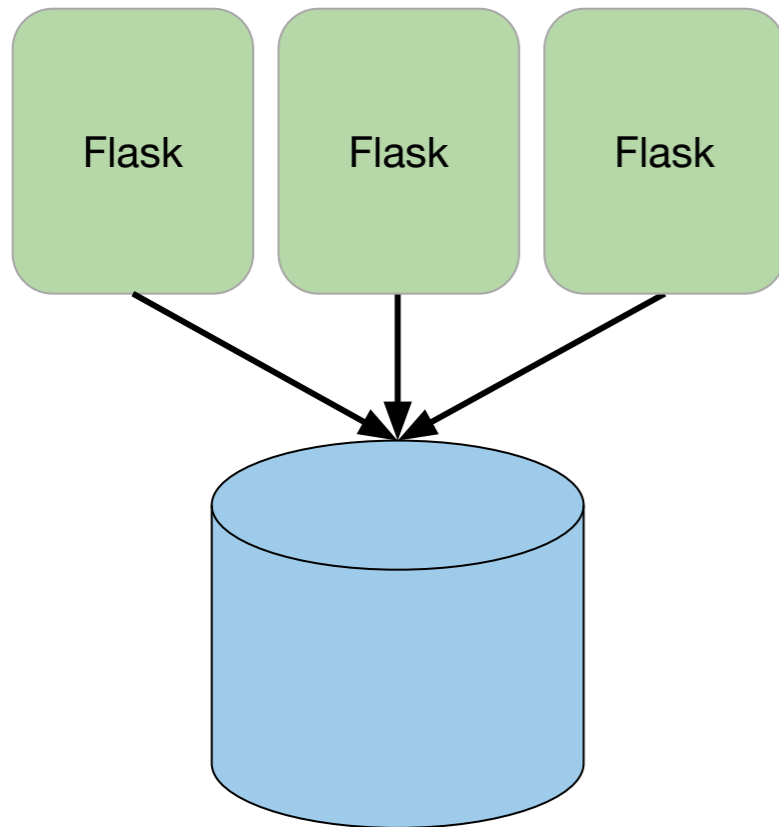
Buy a **faster** computer

Buy **another** computer

# Processing Data in PARALLEL

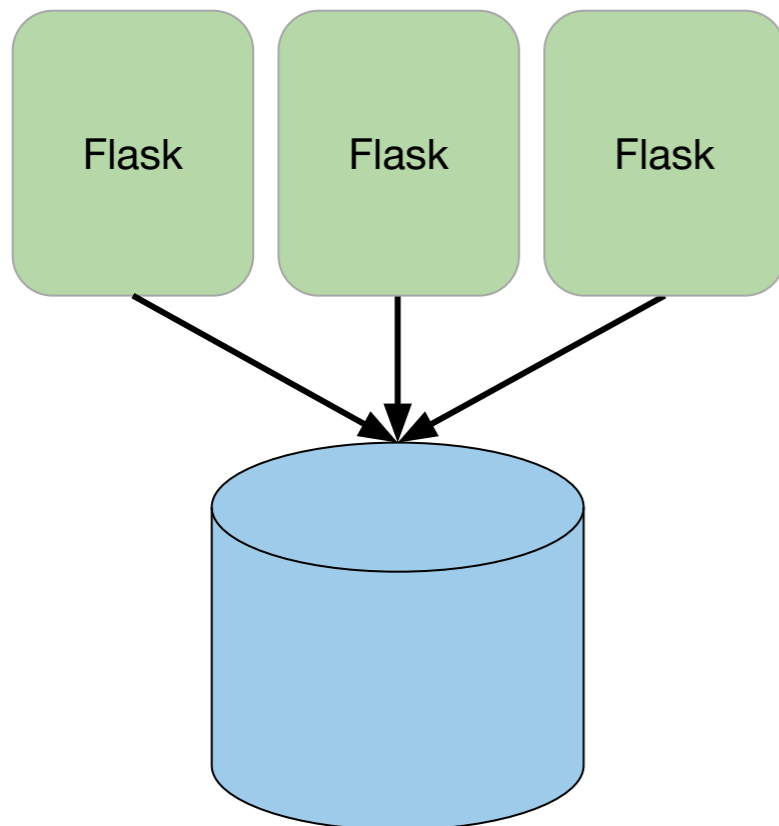


# Processing Data via DB

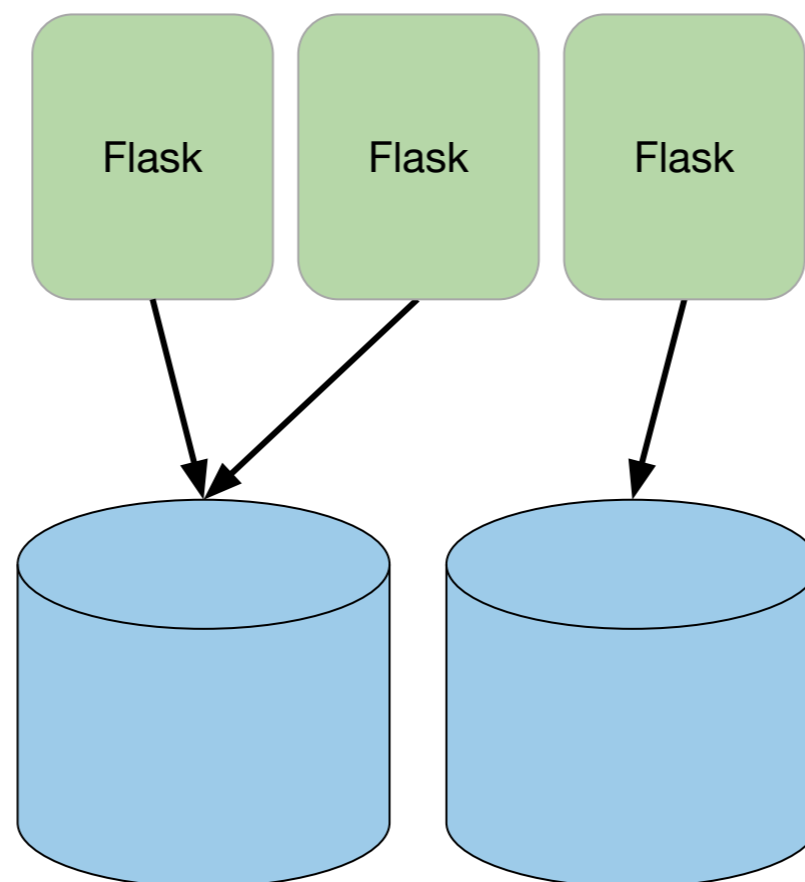


Problems?

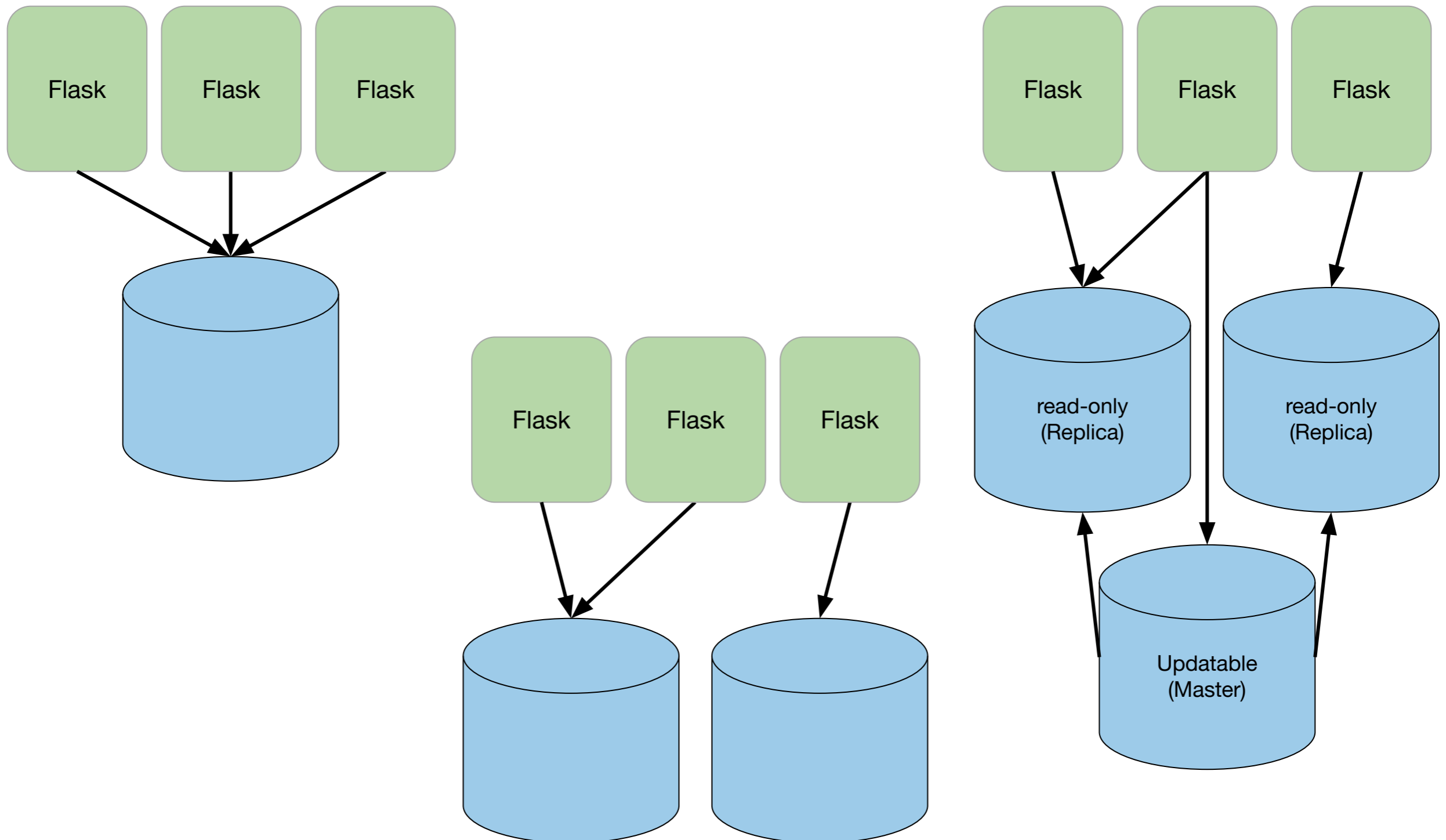
# Processing Data via DB



Problems?



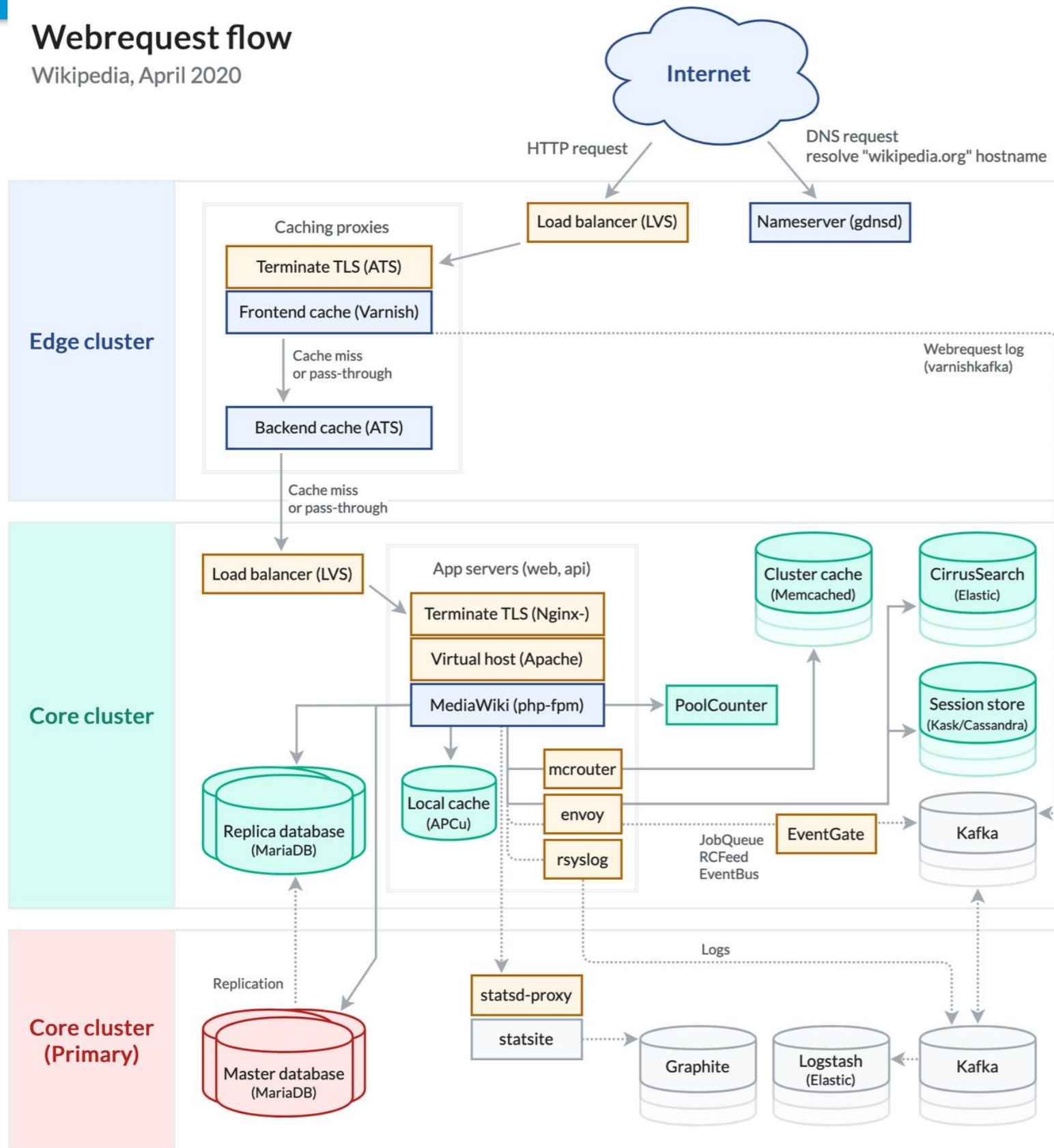
# Processing Data via DB



# How does even Wikipedia work?

## Webrequest flow

Wikipedia, April 2020



[https://meta.wikimedia.org/wiki/Wikimedia\\_servers](https://meta.wikimedia.org/wiki/Wikimedia_servers)  
<https://grafana.wikimedia.org/>