

# SHORT HISTORY of BIG DATA & STREAMING PROGRAMMING TECHNOLOGY



# Jeffrey Ricker

— — —

- 1991BS Mechanical Engineering (Robotics) Tulane University
- 1996US DOD High Performance Computing Modernization Program
- 1997 DARPA Shaolin Project
- 1998 Founded XMLSolutions Corp
- 2004 Founded Distributed Instruments LLC
- 2013 Amazon Big Data
- 2015 Founded The Ricker Lyman Robotic Company

# Objective

— — —

Provide historical context of technology evolution leading up to streaming big data

# Agenda

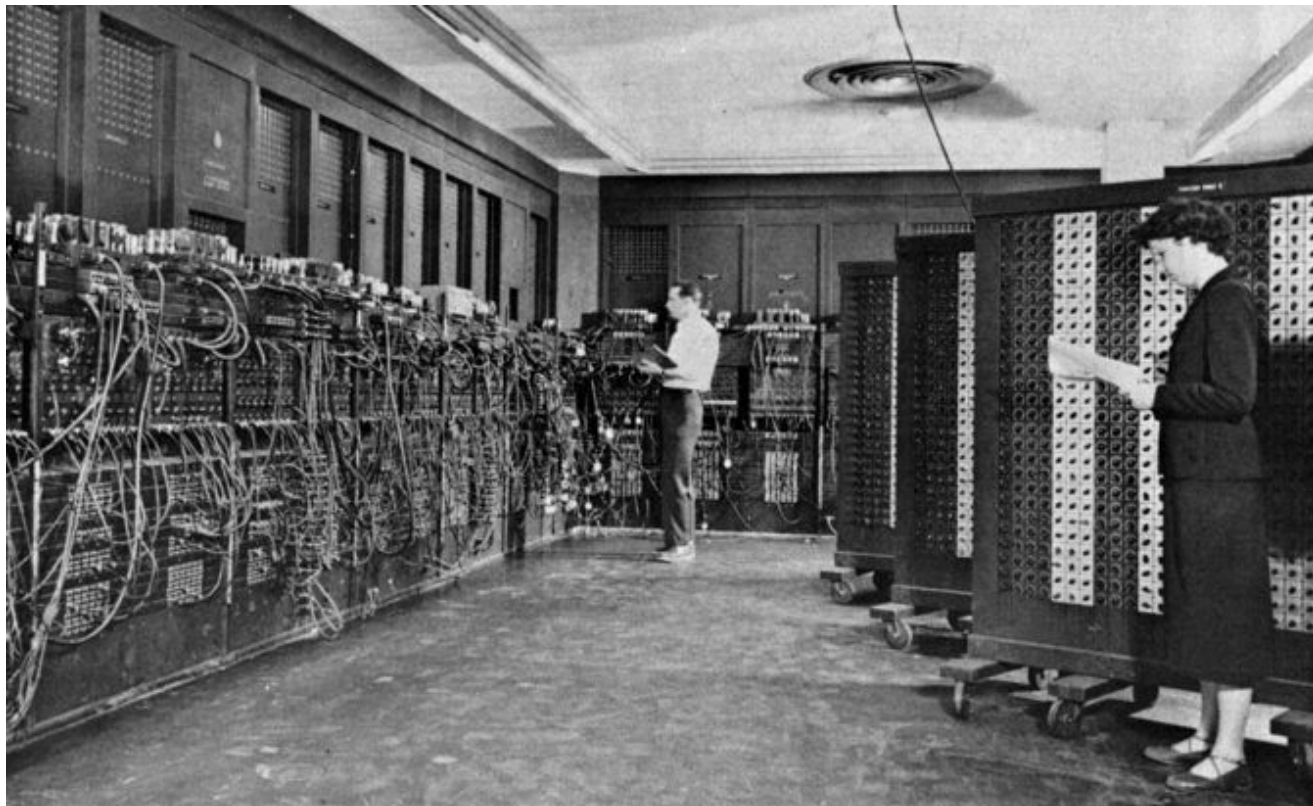
— — —

1. High performance computing
2. Open source
3. Hadoop (big data)
4. Functional programming
5. Streaming programming
6. Why history?

# High Performance Computing

**Size matters**

# 1946 ENIAC



# 1964 CDC 6600



# 1976 Cray 1





# 1991 CM-5



# 1994 Beowulf cluster



# 2018 Summit

Summit has 4,356 nodes, each one equipped with two 22-core **Power9** CPUs, and six **NVIDIA Tesla V100 GPUs**. The nodes are linked together with a Mellanox dual-rail EDR InfiniBand network.



# High Performance Computing

— — —

1940-1970: the first supercomputers

1975-1990: the Cray era

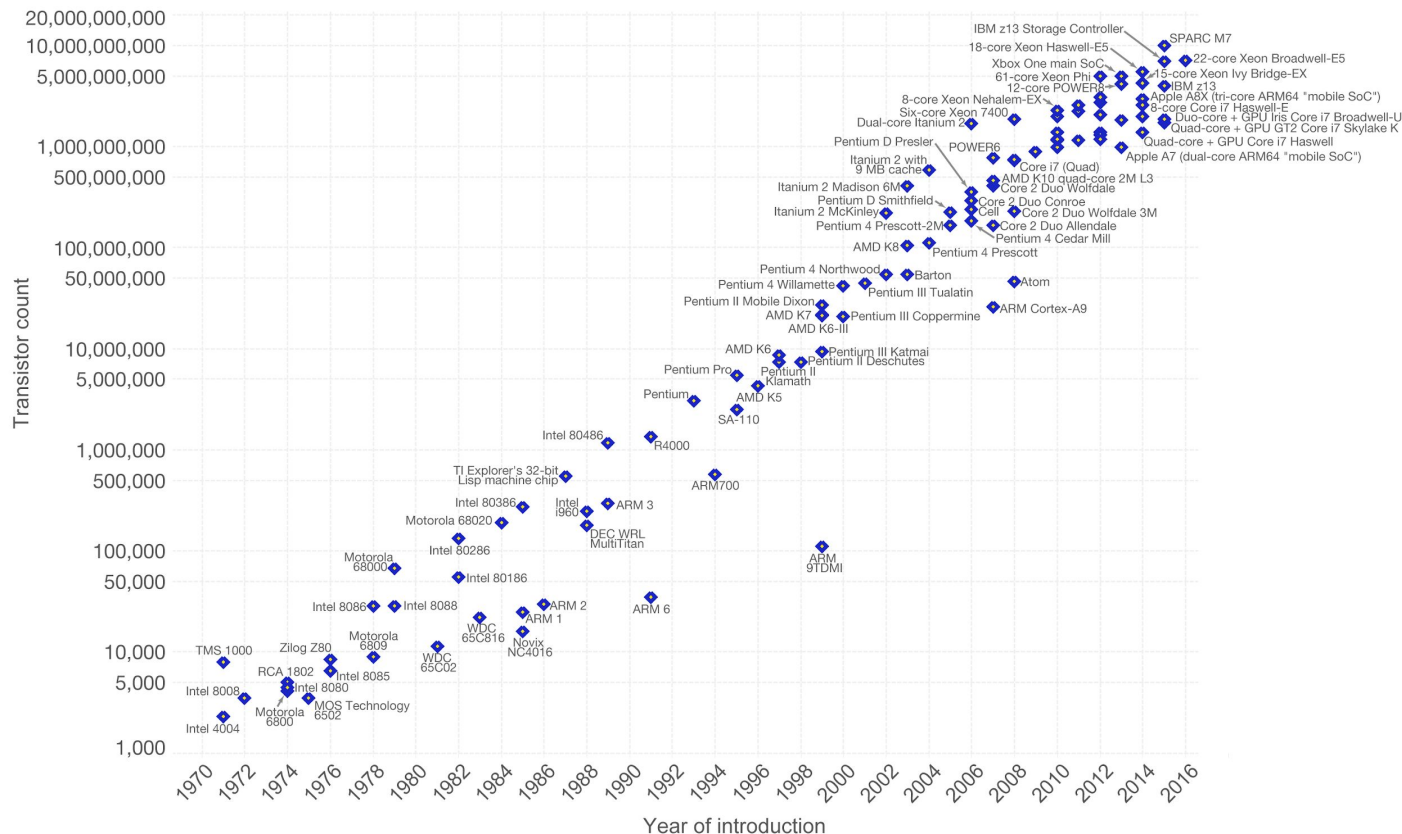
1990-2010: the cluster era

2000-2020: the GPU and hybrid era

2020-: ???

Our World  
in Data

Moore's law describes the empirical regularity that the number of transistors on integrated circuits doubles approximately every two years. This advancement is important as other aspects of technological progress – such as processing speed or the price of electronic products – are strongly linked to Moore's law.

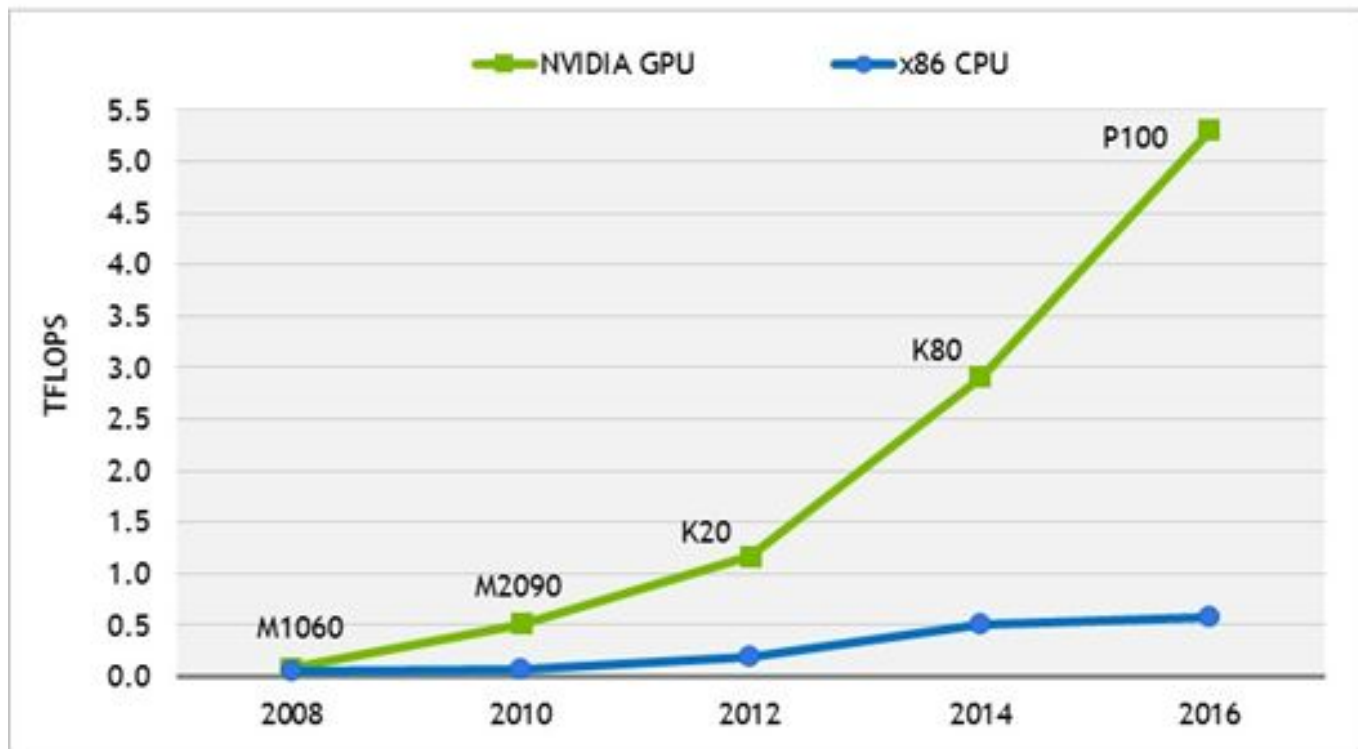


Data source: Wikipedia ([https://en.wikipedia.org/wiki/Transistor\\_count](https://en.wikipedia.org/wiki/Transistor_count))

The data visualization is available at [OurWorldinData.org](https://ourworldindata.org). There you find more visualizations and research on this topic.

Licensed under [CC-BY-SA](#) by the author Max Roser.

# GPU evolution

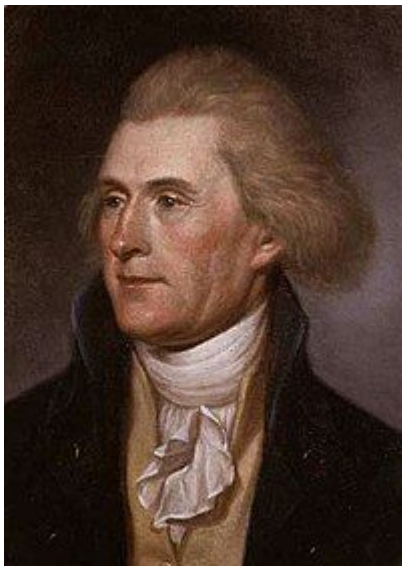


# Open Source

**A business model of innovation**



# 1790



*The United States.*

*To all to whom these Presents shall come. Greeting.*

*Whereas Samuel Hopkins of the City of Philadelphia and State of Pennsylvania hath discovered an Improvement, not known or used before, such Discovery, in the making of Pot ash and Char. ash by a new Apparatus and Process; that is to say, in the making of Char. ash 1<sup>st</sup> by burning the raw Ashes in a Furnace, 2<sup>d</sup> by dispoiling and boiling them when so burnt in Water, 3<sup>d</sup> by drawing off and settling the ley, and 4<sup>th</sup> by boiling the ley into darts which then are the true Char. ash; and also in the making of Pot. ash by fluxing the Char. ash so made as aforesaid, which Operation of burning the raw Ashes in a Furnace, preparatory to their Dispoiling and boiling in Water, is new, leaves little Residuum; and produces a much greater Quantity of salt: These are therefore in pursuance of the Act, entitled "An Act to promote the Progress of useful Arts", to grant to the said Samuel Hopkins, his Heirs, Administrators and Assigns, for the Term of fourteen Years, the sole and exclusive Right and Liberty of using and vending to others the said Discovery, of burning the raw Ashes previous to their being dispoiled and boiled in Water, according to the true Intent and Meaning of the Act aforesaid. In Testimony whereof I have caused these Letters to be made patent, and the Seal of the United States to be hereunto affixed. Given under my Hand at the City of New York this thirty first Day of July in the Year of our Lord one thousand seven hundred & Ninety.*

*G. Washington*

*City of New York July 31<sup>st</sup> 1790.*

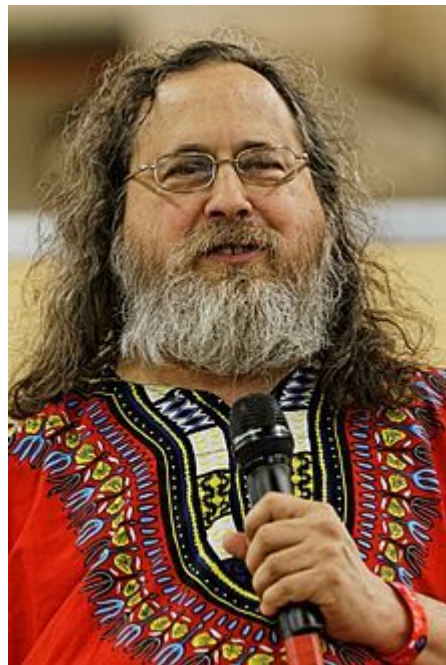
*I do hereby certify that the foregoing Letters patent were delivered to me in pursuance of the Act, entitled "An Act to promote the Progress of useful Arts"; that I have examined the same, and find them conformable to the said Act.*

*Edm: Randolph* Attorney General for the United States.

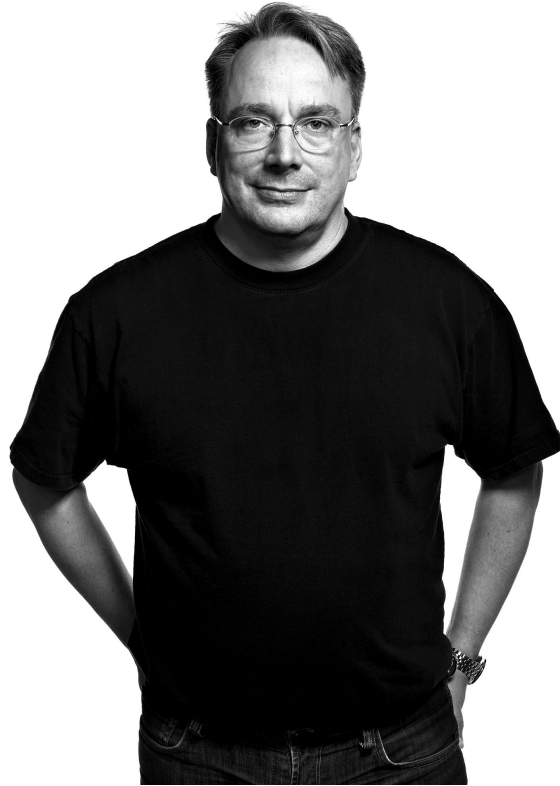
X000001  
July 31, 1790



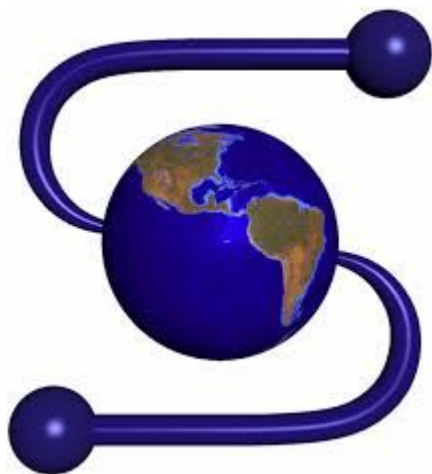
# 1984



1991



1993



1999



# Key events in open source

— — —

1984 Richard Stallman (MIT) starts GNU project

1989 GPL

1991 Linus Torvalds releases Linux

1993 Mosaic browser  
Red Hat founded

1994 Netscape  
MySQL launched

1996 Apache launched

1997 Eric Raymond "The Cathedral and the Bazaar"

1998 Netscape open sources  
Mozilla Firefox

1999 Apache Foundation  
IBM announces \$1 billion investment  
in Linux

**2006 Hadoop incubator**

# Hadoop

**How distributed computing went mainstream**

# Original search



The image is a screenshot of the original Yahoo! homepage. At the top, there is a navigation bar with icons for 'What's New', 'Check Email', 'Personalize', and 'Help'. The 'YAHOO!' logo is prominently displayed in the center. Below the logo, there are several promotional banners: 'Yahoo! Auctions bid & sell for free', 'Family Guy Win 6 days in Hawaii! Go to fox.com FOX', and 'Park Your Domain Free'. A search bar with a 'Search' button and a link to 'advanced search' is located below the banners. A central message reads 'Yahoo! Mail - Get your free e-mail account today!'. Below this, a list of links includes 'Shopping', 'Yellow Pages', 'People Search', 'Maps', 'Travel Agent', 'Classifieds', 'Personals', 'Games', 'Chat', 'Email', 'Calendar', 'Pager', 'My Yahoo!', 'Today's News', 'Sports', 'Weather', 'TV', 'Stock Quotes', and 'more...'. The page is organized into several columns. The left column contains links to 'Arts & Humanities', 'Business & Economy', 'Computers & Internet', 'Education', 'Entertainment', 'Government', and 'Health'. The middle column contains links to 'News & Media', 'Recreation & Sports', 'Reference', 'Regional', 'Science', 'Social Science', and 'Society & Culture'. The right column features a 'In the News' section with a list of headlines and an 'Inside Yahoo!' section with a list of links. The overall design is simple and functional, typical of the early 1990s web.

**YAHOO!**

What's New Check Email Personalize Help

**Yahoo! Auctions**  
bid & sell for free

**Family Guy** Win 6 days in Hawaii! Go to fox.com **FOX**

**Park Your Domain Free**

Search advanced search

**Yahoo! Mail** - Get your free e-mail account today!

[Shopping](#) - [Yellow Pages](#) - [People Search](#) - [Maps](#) - [Travel Agent](#) - [Classifieds](#) - [Personals](#) - [Games](#) - [Chat](#) - [Email](#) - [Calendar](#) - [Pager](#) - [My Yahoo!](#) - [Today's News](#) - [Sports](#) - [Weather](#) - [TV](#) - [Stock Quotes](#) - [more...](#)

**Arts & Humanities**  
[Literature](#), [Photography](#)...

**Business & Economy**  
[Companies](#), [Finance](#), [Jobs](#)...

**Computers & Internet**  
[Internet](#), [WWW](#), [Software](#), [Games](#)...

**Education**  
[Universities](#), [K-12](#), [College Entrance](#)...

**Entertainment**  
[Cool Links](#), [Movies](#), [Humor](#), [Music](#)...

**Government**  
[Military](#), [Politics](#), [Law](#), [Taxes](#)...

**Health**  
[Medicine](#), [Diseases](#), [Drugs](#), [Fitness](#)...

**News & Media**  
[Full Coverage](#), [Newspapers](#), [TV](#)...

**Recreation & Sports**  
[Sports](#), [Travel](#), [Autos](#), [Outdoors](#)...

**Reference**  
[Libraries](#), [Dictionaries](#), [Quotations](#)...

**Regional**  
[Countries](#), [Regions](#), [US States](#)...

**Science**  
[Biology](#), [Astronomy](#), [Engineering](#)...

**Social Science**  
[Archaeology](#), [Economics](#), [Languages](#)...

**Society & Culture**  
[People](#), [Environment](#), [Religion](#)...

**In the News**

- [King Hussein of Jordan dies](#)
- [Online: Lewinsky video testimony](#)
- [NASA comet mission](#)
- [NBA season opens](#)
- [Weekend's top movies](#)

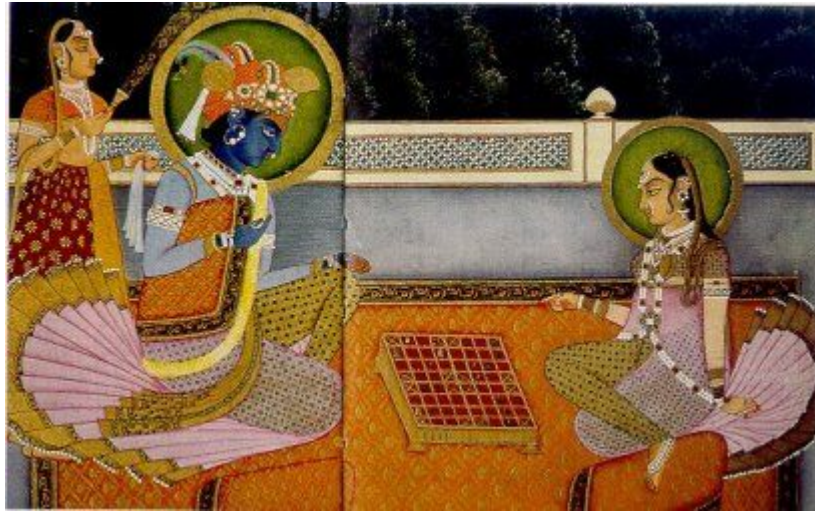
[more...](#)









**Inside Yahoo!**

- [Y! Personals](#) - find a Valentine
- [Shop](#) for your Valentine
- [Y! Clubs](#) - create your own

[more...](#)

# Legend of Paal Paysam



								128
256	512	1,024	2,048	4,096	8,192	16,384	32,768	
64K	128K	256K	512K	1M	2M	4M	8M	
16M	32M	64M	128M	256M	512M	1G	2G	
4G	8G	16G	32G	64G	128G	256G	512G	



# 1997 Lucene

— — —

Doug Cutting builds full text search library

Analyze ordinary text with the purpose of building an index.

Index is a data structure that maps each term to its location in text, so that when you search for a term, it immediately knows all the places where that term occurs.

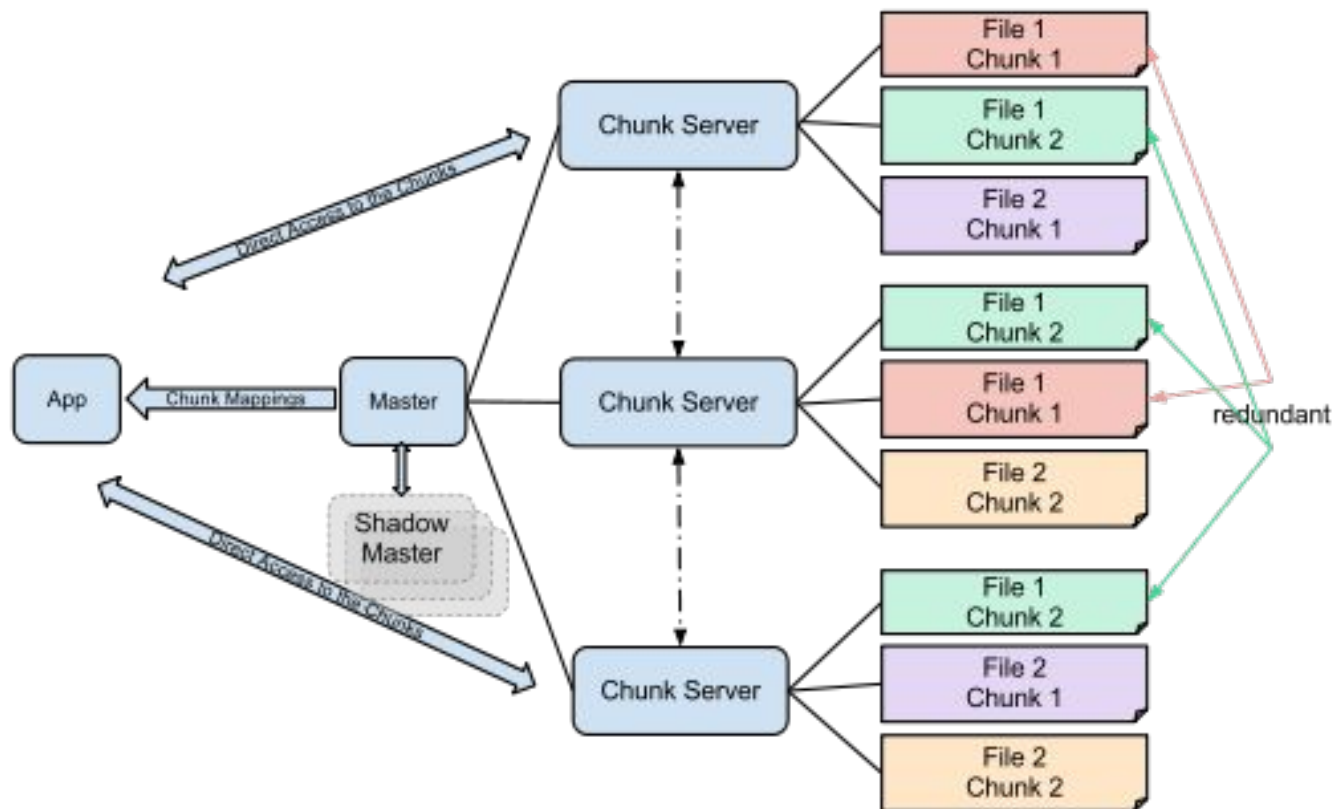
Add Nutch: a web crawler

# 2003 Google File System

— — —

- **schemaless** with no predefined structure, i.e. no rigid schema with tables and columns (and column types and sizes)
- **durable** once data is written it should never be lost
- capable of handling **component failure** without human intervention (e.g. CPU, disk, memory, network, power supply, MB)
- **automatically rebalanced** to even out disk space consumption throughout cluster

# Google file system



# 2004 Google MapReduce

— — —

The three main problems that the MapReduce paper solved are:

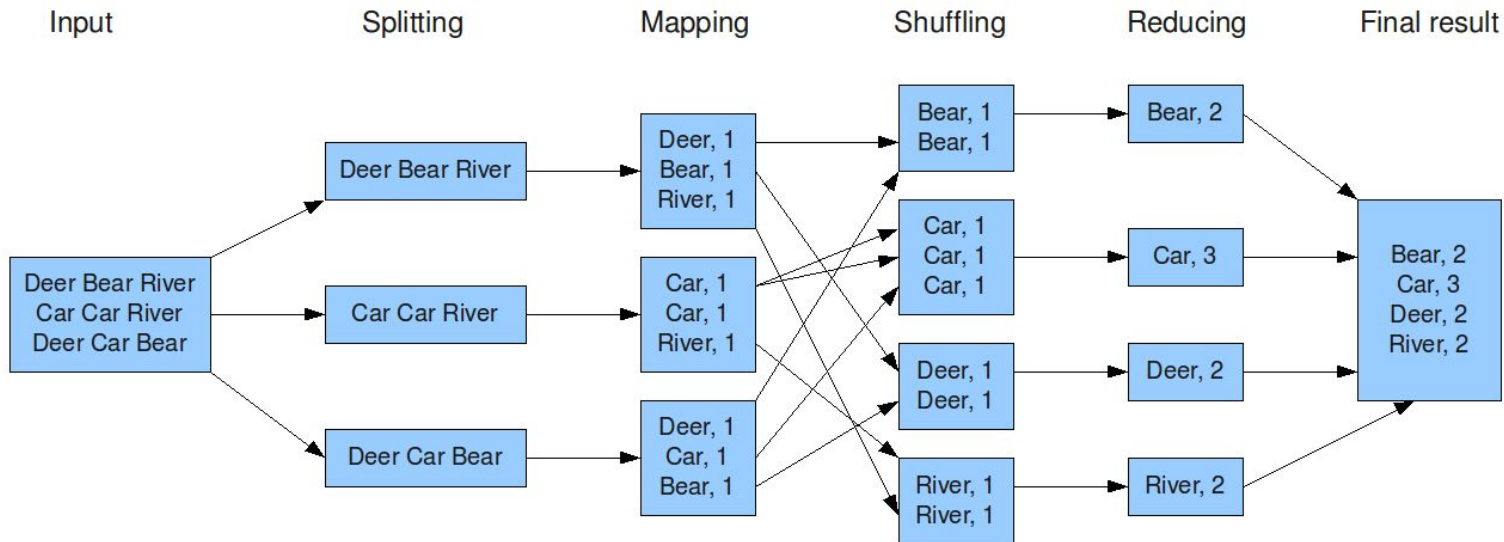
1. Parallelization – how to parallelize the computation
2. Distribution – how to distribute the data
3. Fault-tolerance – how to handle component failure

A program is sent to where the data resides.

There's simply too much data to be moved around.

# MapReduce

The overall MapReduce word count process



# Pre Hadoop

— — —

- 1997 Lucene started by Doug Cutting
- 2000 Lucene released to Source Forge
- 2001 Lucene becomes Apache project <http://lucene.apache.org/>
- 2001 Cutting and Mike Cafarella start Apache Nutch
- 2003.10 Google File System published
- 2004 Nutch Distributed File System
- 2004.01 Scala released
- 2004.12 Map Reduce published <https://ai.google/research/pubs/pub62>
- 2005.06 MapReduce integrated into Nutch

# Post Hadoop

— — —

2006.02 Hadoop incubator released

2006 Bigtable published <https://ai.google/research/pubs/pub27898>

2007.02 Yahoo reports 1000 node Hadoop cluster

2008.01 Hadoop becomes a top level Apache project

2008 HBase joins Hadoop

2008.05 ZooKeeper

2008.10 Pig (from Yahoo) and Hive (from Facebook)

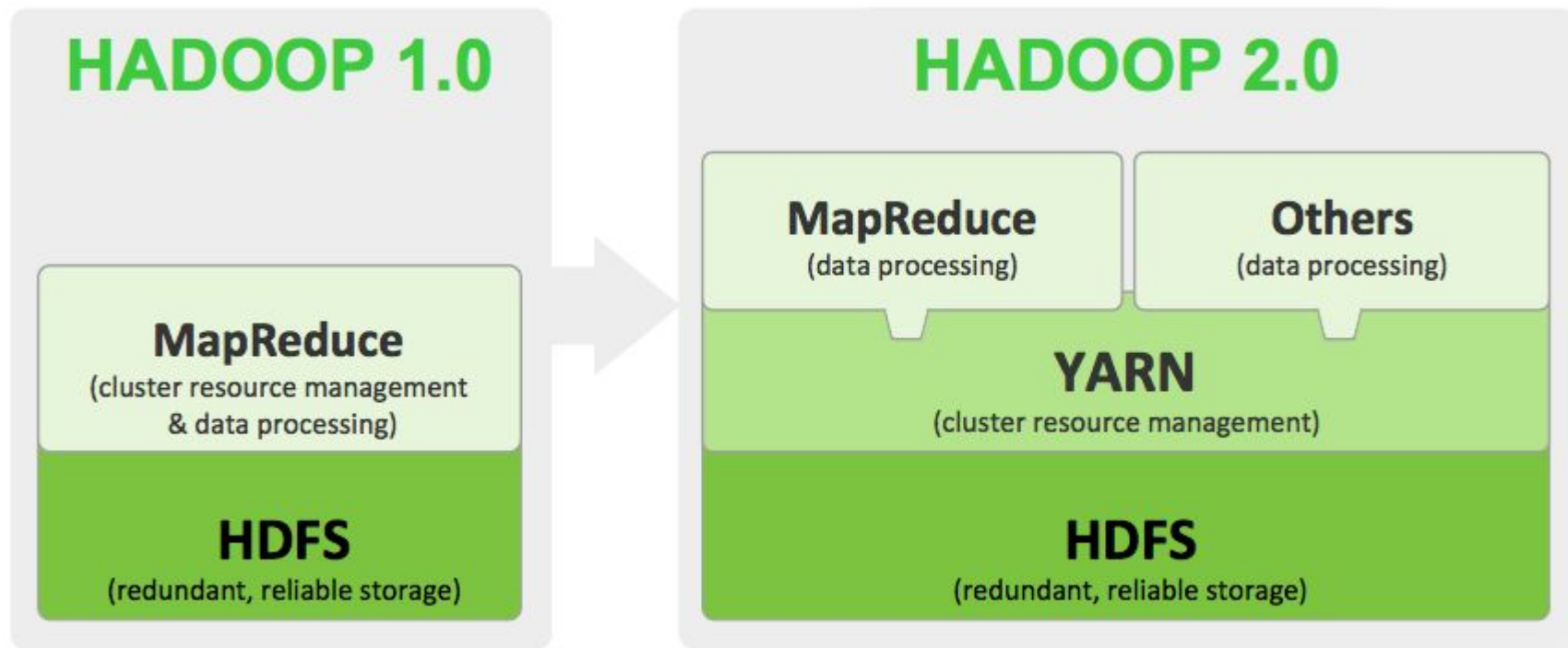
2008 Cloudera founded

2012 Yahoo Hadoop cluster reaches 42,000 nodes

2012.08 YARN becomes subproject

# 2012 YARN

— — —





# Functional programming

What is old is new again

# Programming paradigms

— — —

- Imperative
- Structured
- Procedural
- Object-oriented
- Event-driven
- Declarative
- Functional
- Reactive

# Gottfried von Leibniz

— — —

(1) Create a 'universal language' in which all possible problems can be stated.

(2) Find a decision method to solve all the problems stated in the universal language.

*Entscheidungsproblem*



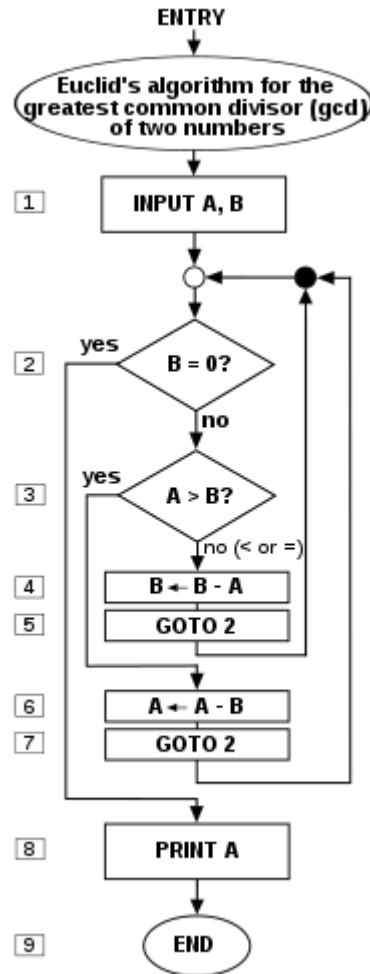
# 1936 Church & Turing



# Turning machine



# Algorithms



# Lambda calculus

— — —

- A formal system in mathematical logic for expressing computation based on function abstraction and application using variable binding and substitution.
- The smallest universal programming language of the world. The  $\lambda$  calculus consists of a single transformation rule (variable substitution) and a single function definition scheme.

# Lambda examples

— — —

first of (sort (append ('dog', 'rabbit') (sort (('mouse', 'cat')))))

→ first of (sort (append ('dog', 'rabbit') ('cat', 'mouse')))

→ first of (sort ('dog', 'rabbit', 'cat', 'mouse'))

→ first of ('cat', 'dog', 'mouse', 'rabbit')

→ 'cat'.

(7 + 4) \* (8 + 5 \* 3)

→ 11 \* (8 + 5 \* 3)

→ 11 \* (8 + 15)

→ 11 \* 23

→ 253.



# Lisp

— — —

```
(defun factorial (n)
  (if (= n 0) 1
      (* n (factorial (- n 1)))))
```

# Functional programming languages

— — —

1932 Lambda calculus -- Alonzo Church

1958 LISP -- John McCarthy

1970 Scheme

1986 Erlang

1990 Haskell

1995 JavaScript

**2004 Scala**

2005 F#

2007 Clojure

# Stream Programming

**What, when, why, how?**

# A rose by any other name...

— — —

- Stream programming
- Stream processing
- Real-time analytics
- Streaming analytics
- Complex event processing (CEP)
- Real-time streaming analytics
- Event processing

# Earlier CEP frameworks

— — —

- 2002 Aurora
- 2005 Borealis
- 2005 Apama
- 2007 Cayuga
- 2008 Esper
- 2011 Apache S4 (Yahoo)

# What is CEP?

— — —

## Trivial

React to a button pushed

0001 Push

0005 Push

## Complex

React to a button pushed 3 times in 10 seconds

0012 Push

0014 Push

ALERT

0021 Push

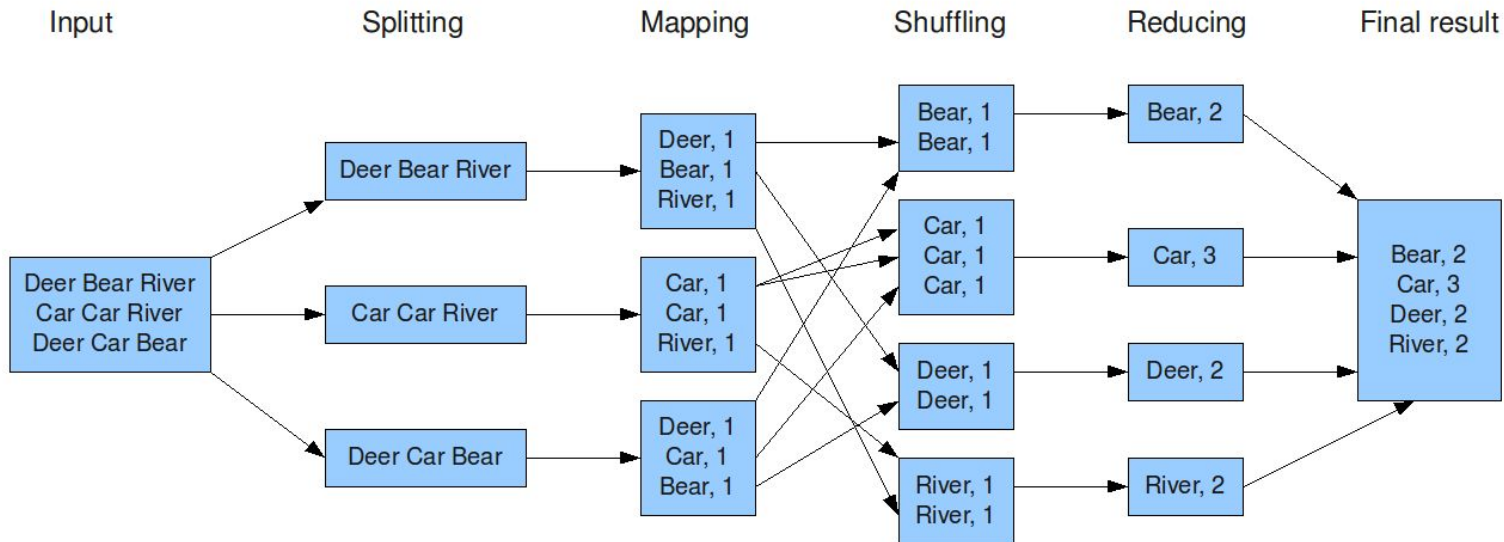
alert again?

0033 Push

# MapReduce

— — —

The overall MapReduce word count process



# Workflows

— — —

Oozie

Luigi

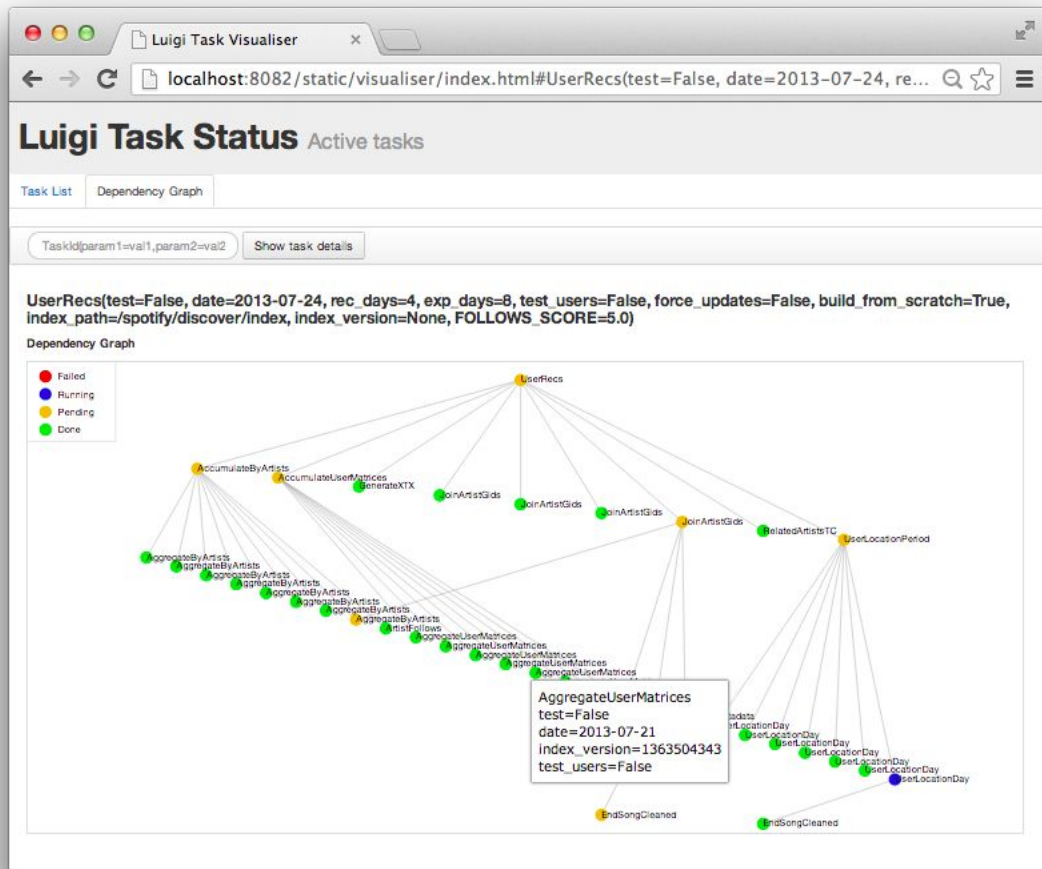
Azkaban

Airflow

Pinball

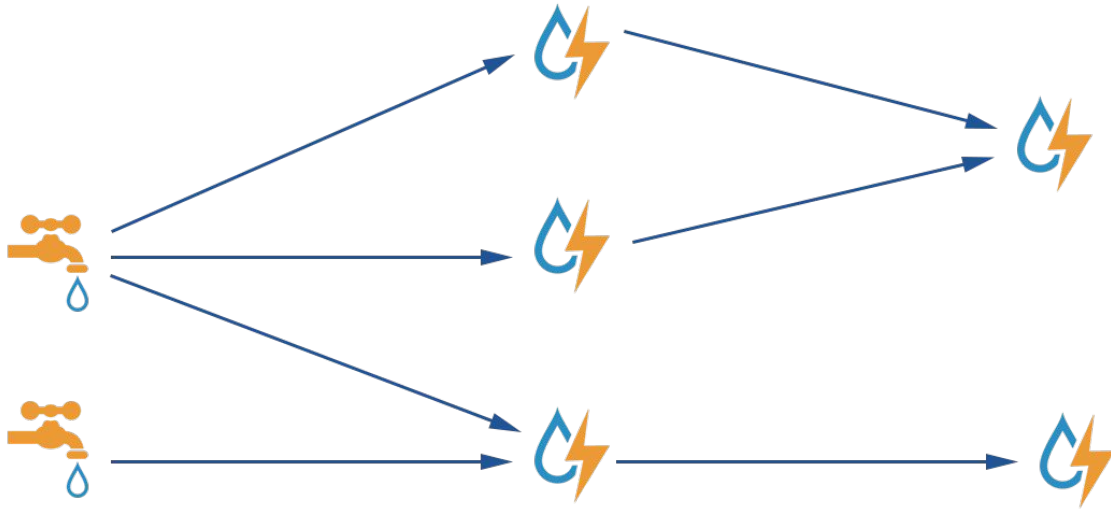
Cascading

Taskflow





# 2011 Apache Storm



# Streaming frameworks

— — —

2011 Apache Storm

2014 **Apache Spark**

Apache Samza

2015 Apache Flink

Apache Nifi

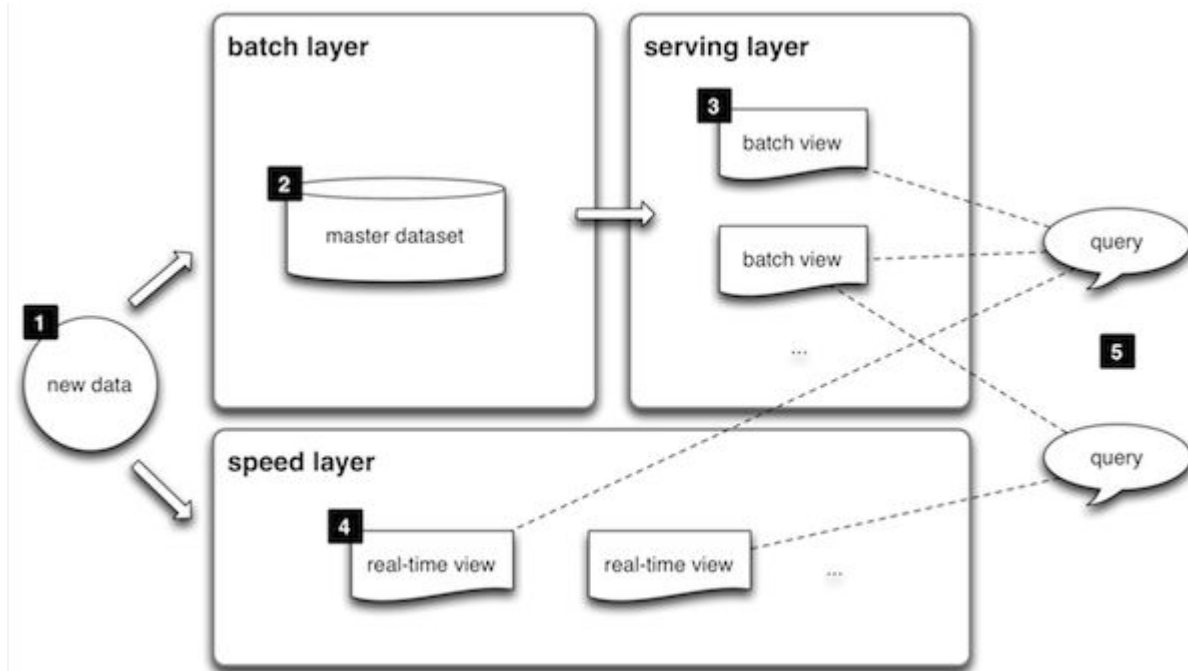
2016 Apache Gearpump

Apache Apex

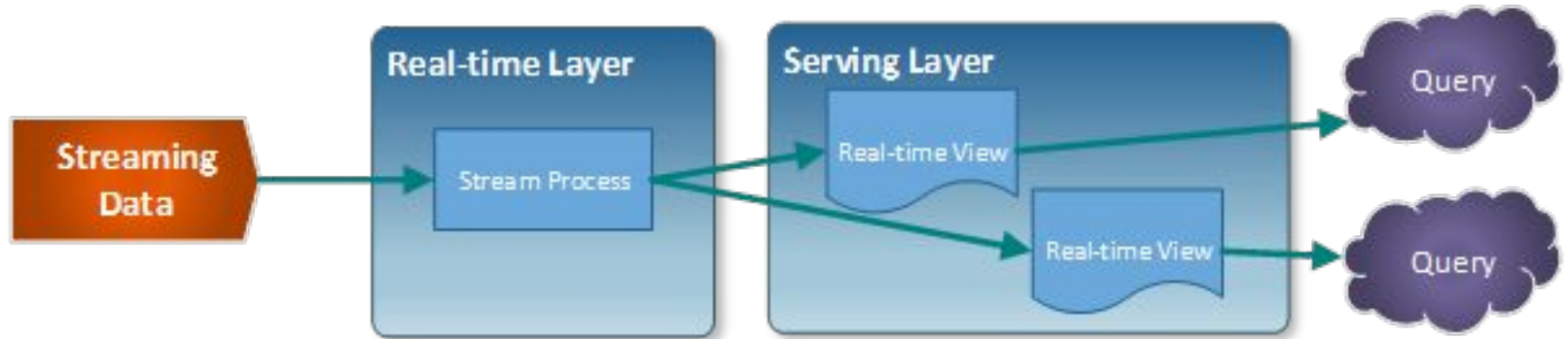
**Kafka Streams**

Akka Streams

# 2013 Lambda Architecture



# 2015 Kappa architecture



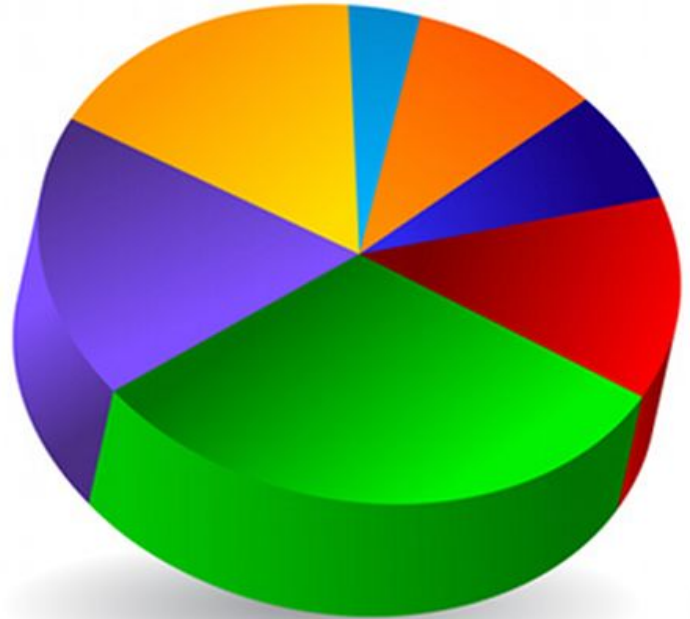
# Why history?

**The only constant is change**

# Pie metaphor

— — —

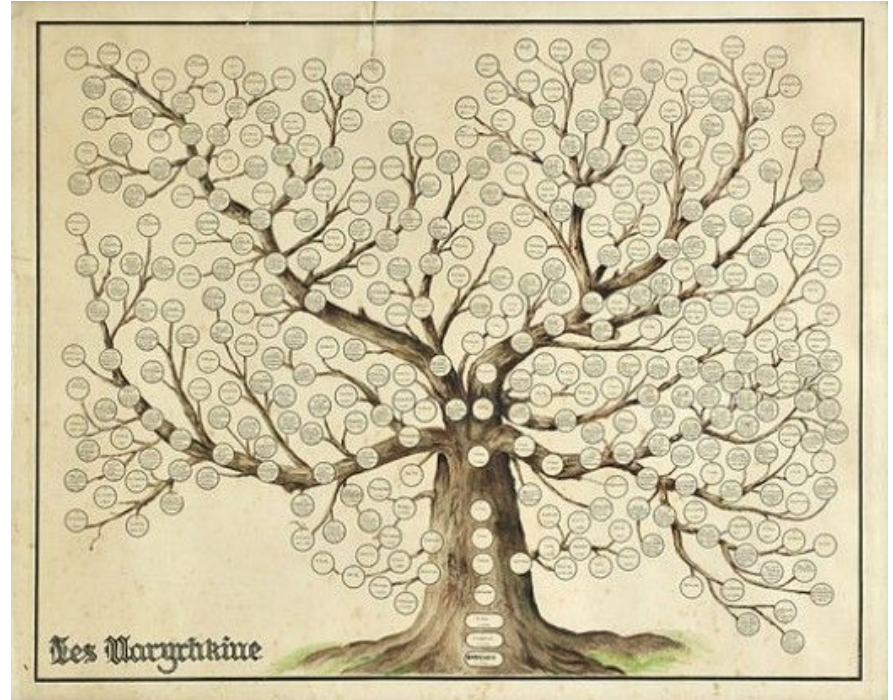
- Different technologies fill different roles, like slices of a pie
- Where slices meet neither solution is a perfect match
- The pie is always expanding
- The boundary between slices become large gaps
- New technologies arise to fill the gaps



# Tree metaphor

— — —

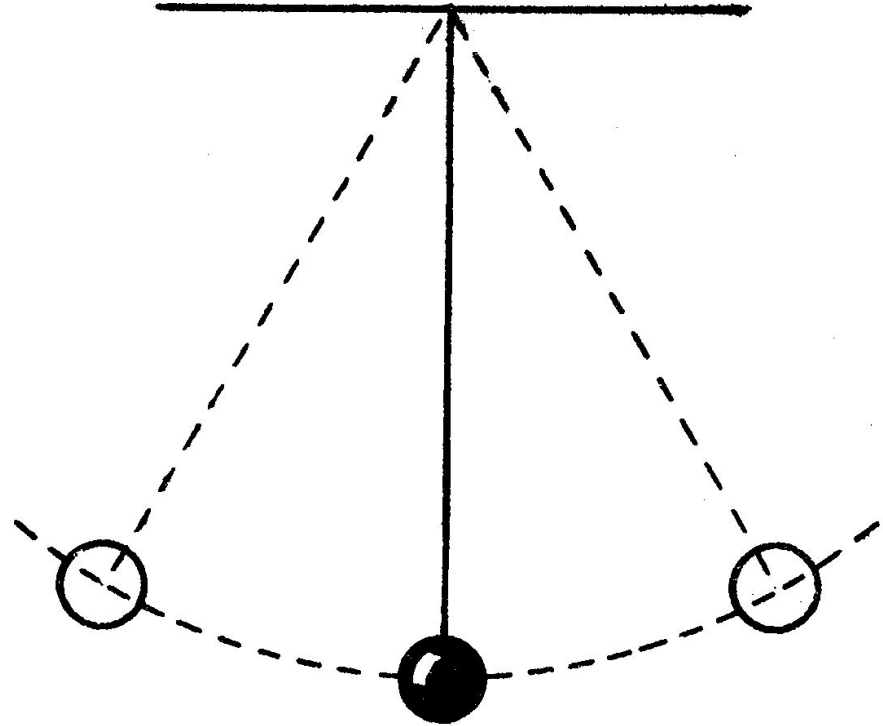
- There are many, many instances of the same thing
- Focus on root or branch technology before dealing with leaves
- What is old is new again



# Pendulum metaphor

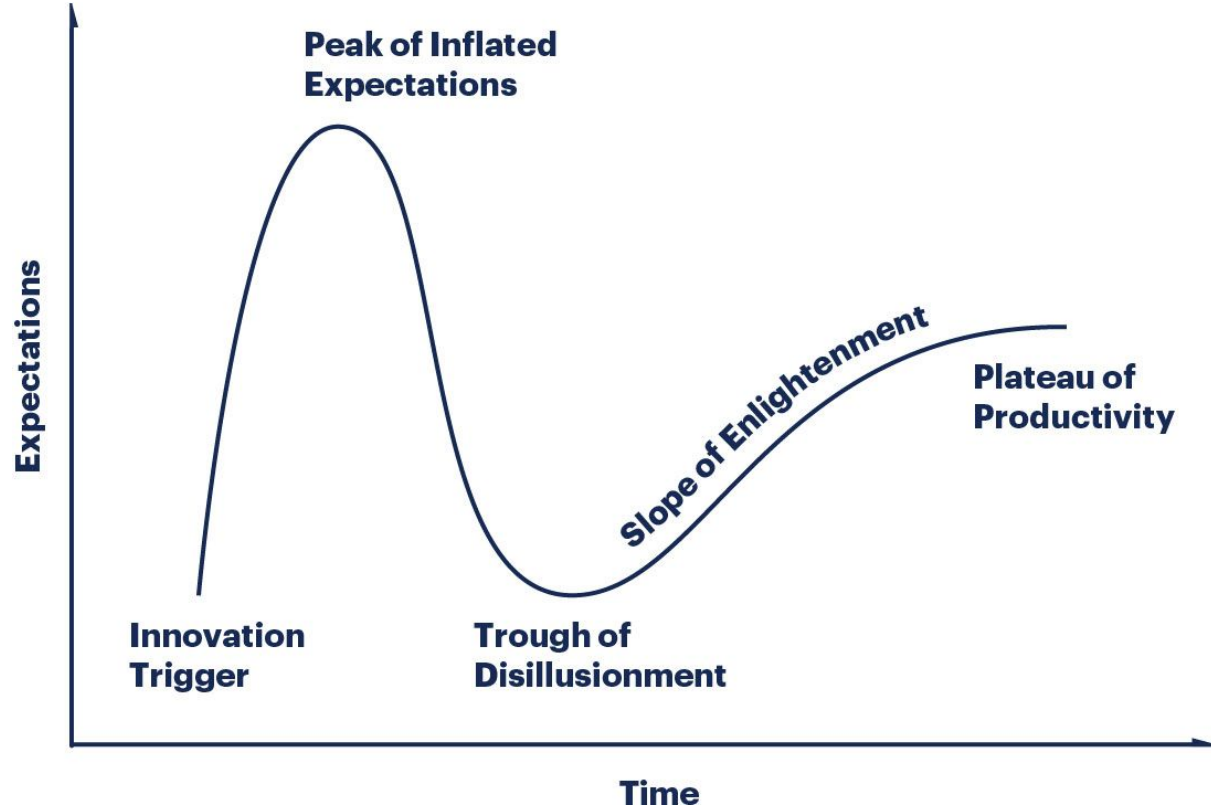
---

- Markets oscillate
- Centralized to decentralized
- Structured to unstructured
- Consolidated to fractured

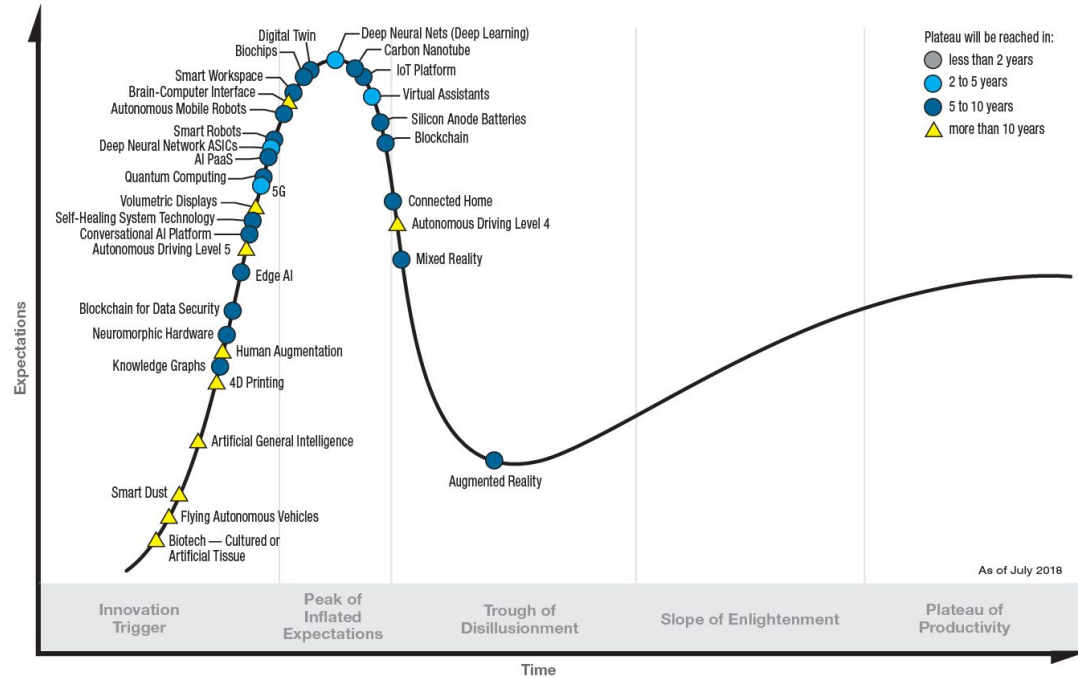




# The hype cycle



## Hype Cycle for Emerging Technologies, 2018

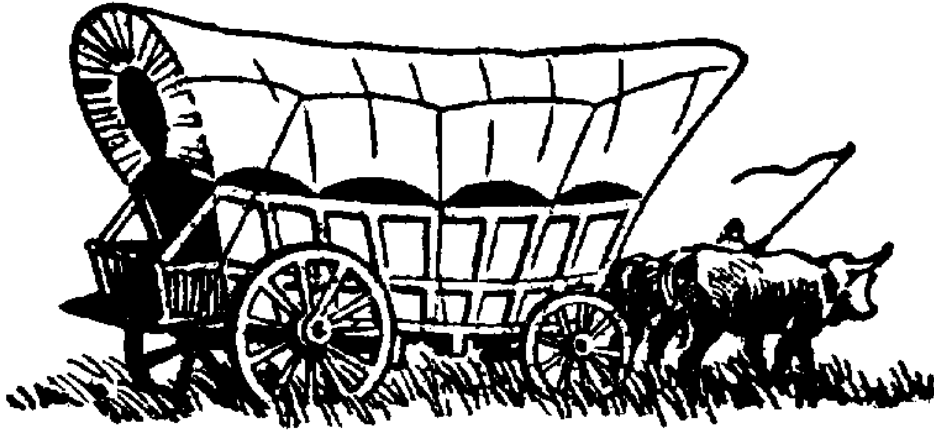


[gartner.com/SmarterWithGartner](https://gartner.com/SmarterWithGartner)

Source: Gartner (August 2018)  
© 2018 Gartner, Inc. and/or its affiliates. All rights reserved.

**Gartner.**

# Build wagons, not cabins



***THANK  
YOU!***

