



# Innovation & Research Infrastructures

## A Google perspective

Ed Parsons

Geospatial Technologist, Google Research, London

[eparsons@google.com](mailto:eparsons@google.com)



+ed parsons



@edparsons

Google

If I have seen further it is by  
standing on the shoulders of  
infrastructures.



# 1. Hybrid Research at Google

- The goal of research at Google is to bring **significant, practical** benefits to our **users**, and to do so rapidly, within a few years at most.
- Approach is iterative and usually involves producing near-production code from day one.
- Single team takes product from Research to Production.
- Requires computing and storage infrastructure at operational scale

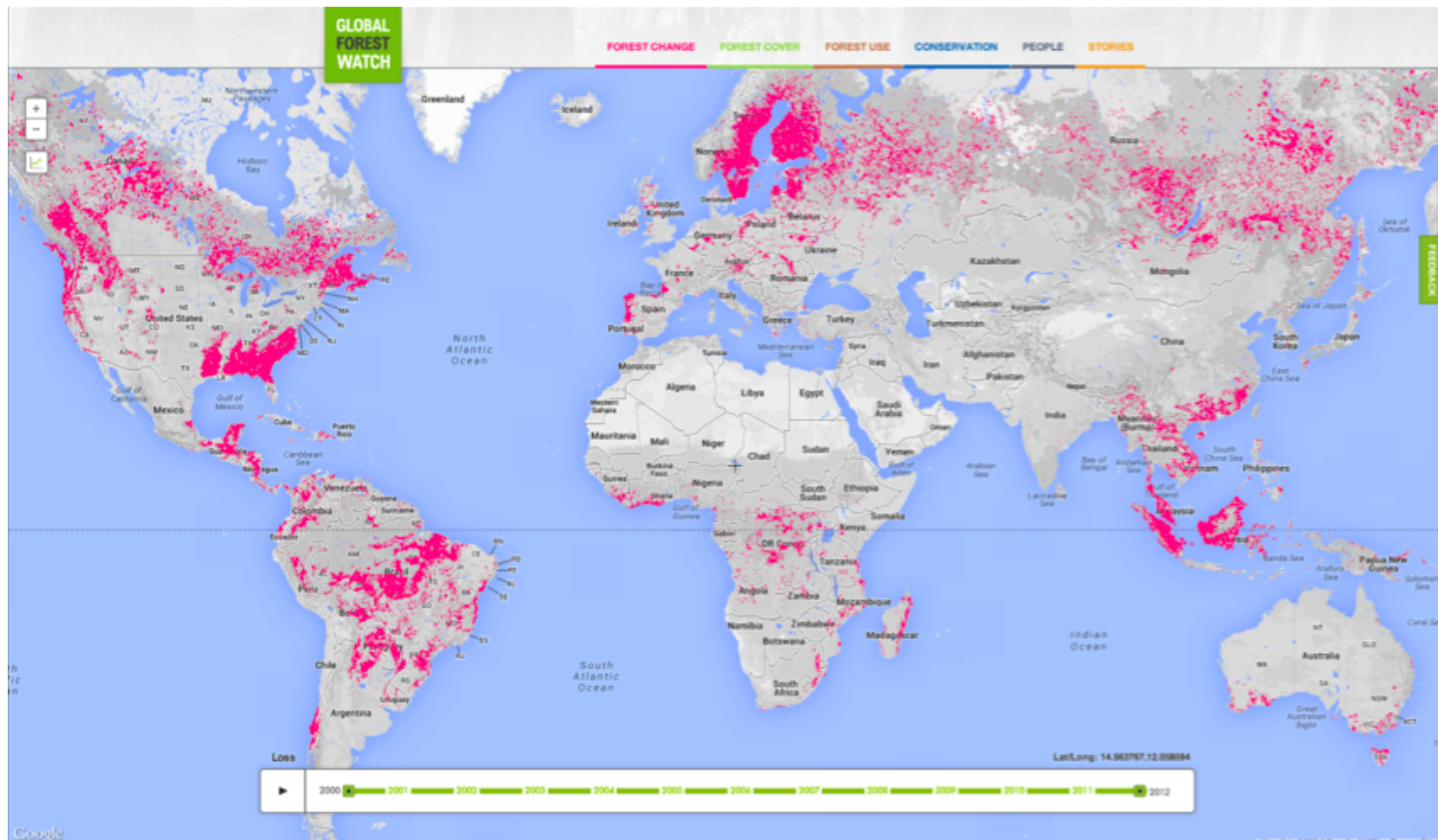


## 2. Research Infrastructures that Scale..

# Open, Accessible Platforms

Technology & Cultural





Small teams



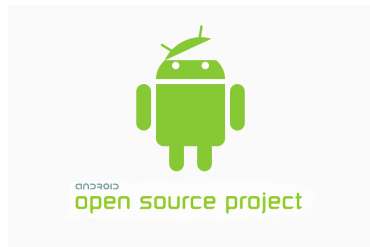
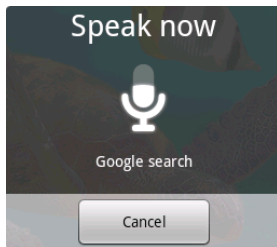
Big Data

### 3. Measure success..

- Academic, Commercial or ideally both !
- Open Source Projects, Industry standards development



Jeffrey Dean and Sanjay Ghemawat  
jeff@google.com, sanjay@google.com  
Google, Inc.



**Abstract**

MapReduce is a programming model and an associated implementation for processing and generating large data sets. Users specify a *map* function that processes a key/value pair to generate a set of intermediate key/value pairs, and a *reduce* function that merges all intermediate values associated with the same intermediate key. Many real world tasks are expressible in this model, as shown in the paper.

Programs written in this functional style are automatically compiled and executed on a large cluster of commodity hardware. The system takes care of the details of scheduling the program, scheduling the program on the machines, handling the program's required inter-machine communication, and managing the program's execution on the hardware system.

Programs written in this functional style are automatically compiled and executed on a large cluster of commodity hardware. The system takes care of the details of scheduling the program, scheduling the program on the machines, handling the program's required inter-machine communication, and managing the program's execution on the hardware system.

As a reaction to this complexity, we designed a new abstraction that allows us to express the simple computations we were trying to perform but hides the messy details of parallelization, fault-tolerance, data distribution and load balancing in a library. Our abstraction is inspired by the *map* and *reduce* primitives present in Lisp and many other functional languages. We realized that most of our computations involved applying a *map* operation to each logical "record" in our input in order to compute a set of intermediate key/value pairs, and then applying a *reduce* operation to all the values that shared the same key, in order to combine the derived data ap-



Ed Parsons

Works at Google  
Attended Cranfield University  
Lives in Teddington

2,576 have him in circles



[www.google.com/+EdParsons](http://www.google.com/+EdParsons)

[eparsons@google.com](mailto:eparsons@google.com)

[@edparsons](https://twitter.com/edparsons)