

# Data Science Research Infrastructure

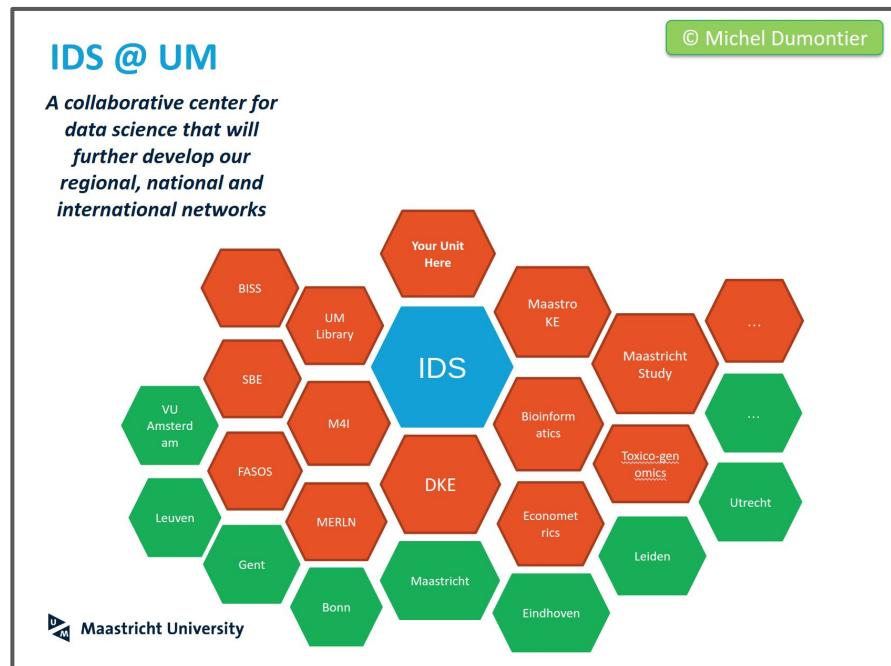
University Meeting Eindhoven - July 2<sup>nd</sup> 2019  
Alexander Malic - Institute Of Data Science @ Maastricht University

# Partners



# Introduction

- Self taught university drop out
- 18 years of professional experience (8 years consulting, 10 years in Research-IT at multinational company)
- Simple solutions for complex problems
- Automate repetitive tasks both for myself and my customers (internal/external)
- Shared vision with Michel



# How does Google do it?

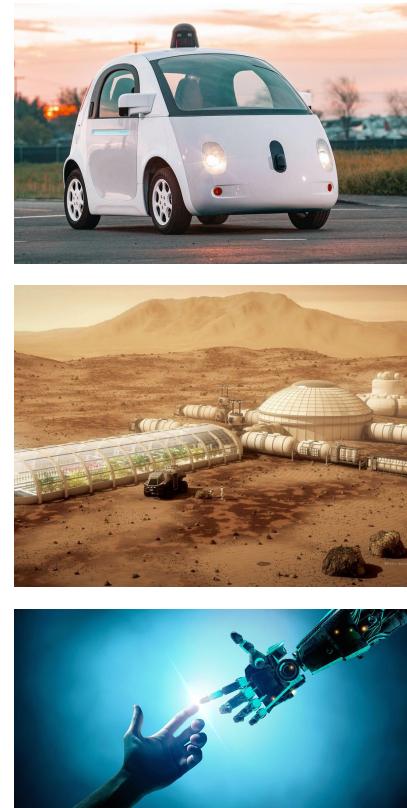
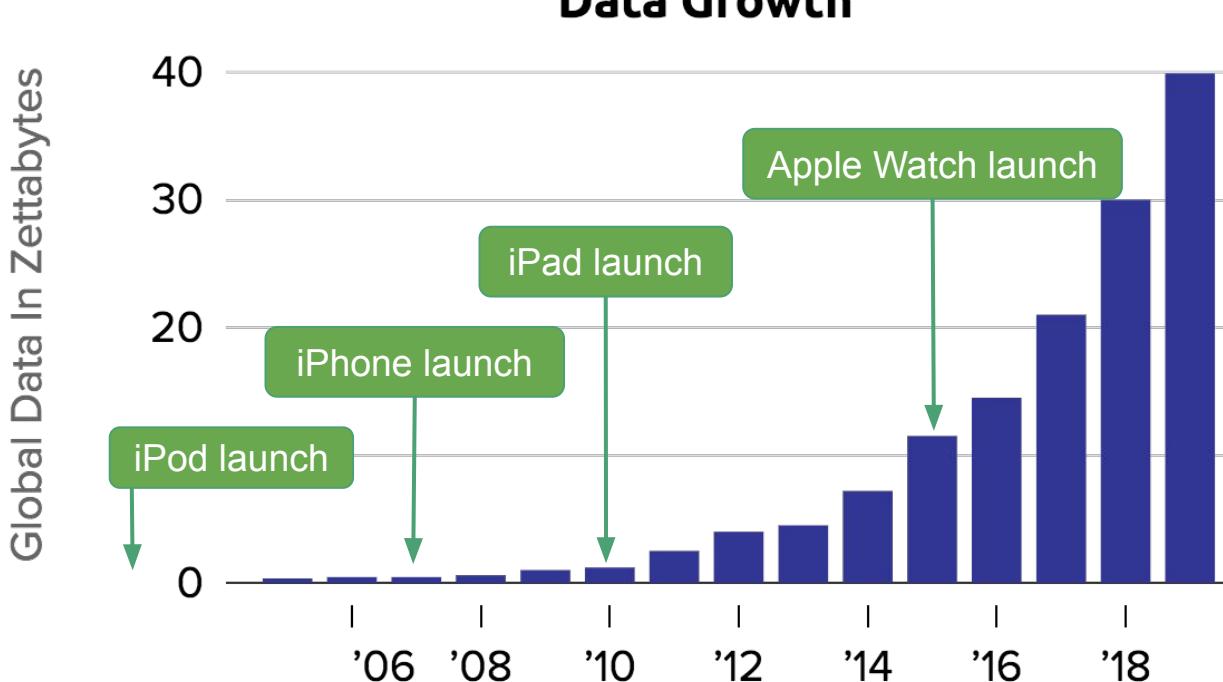
- Scalable & distributed File-System:
  - The Google File system <sup>1</sup>
  - Apache Hadoop
- Scalable & distributed data-base:
  - Google BigTable <sup>2</sup>
  - Apache HBase database
- Scalable computing:
  - Google Borg <sup>3</sup>
  - Docker/Kubernetes

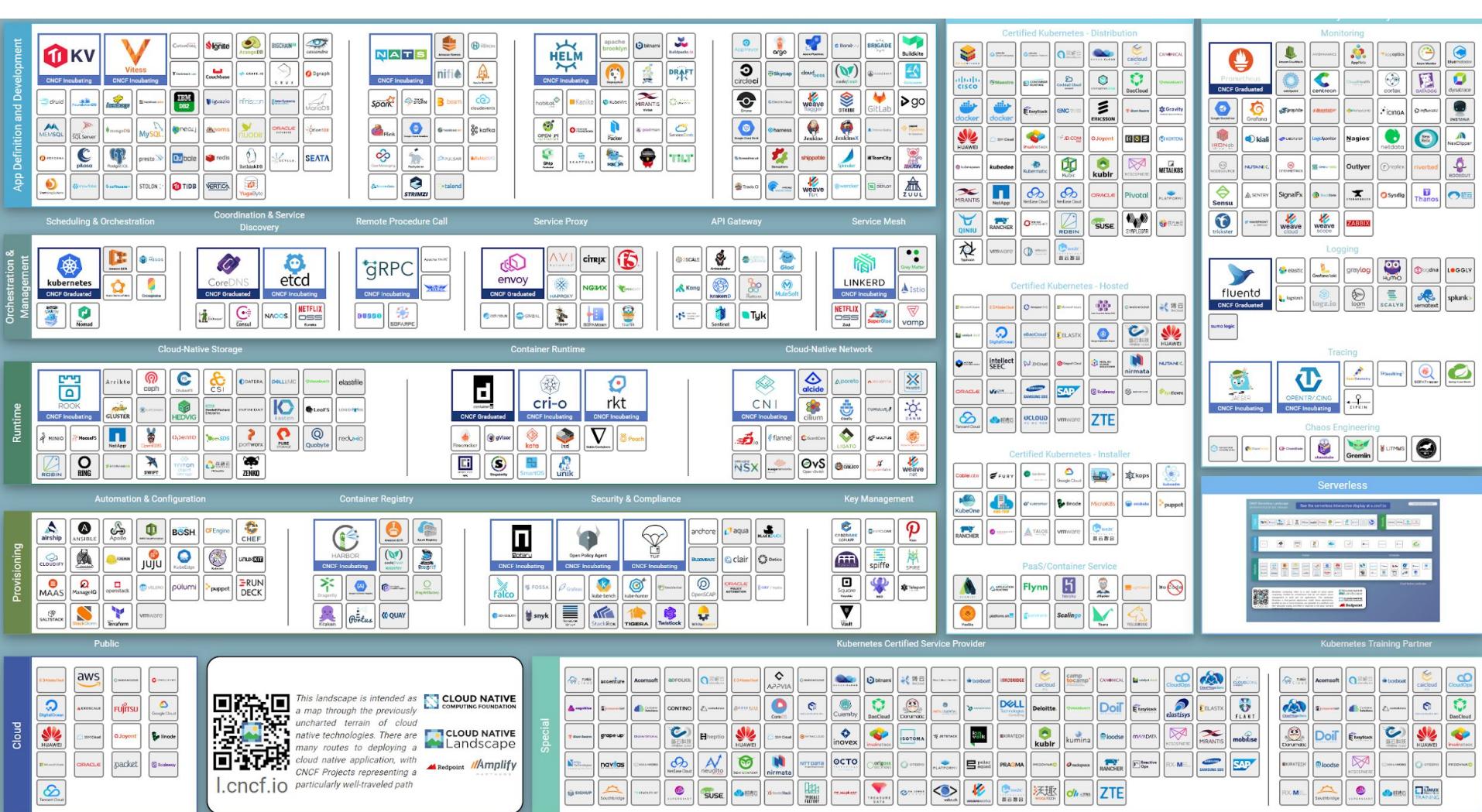
<sup>1</sup> The Google File System - <https://ai.google/research/pubs/pub51>

<sup>2</sup> Bigtable: A Distributed Storage System for Structured Data - <https://ai.google/research/pubs/pub27898>

<sup>3</sup> Large-scale cluster management at Google with Borg - <https://ai.google/research/pubs/pub43438>

# Data growth problem





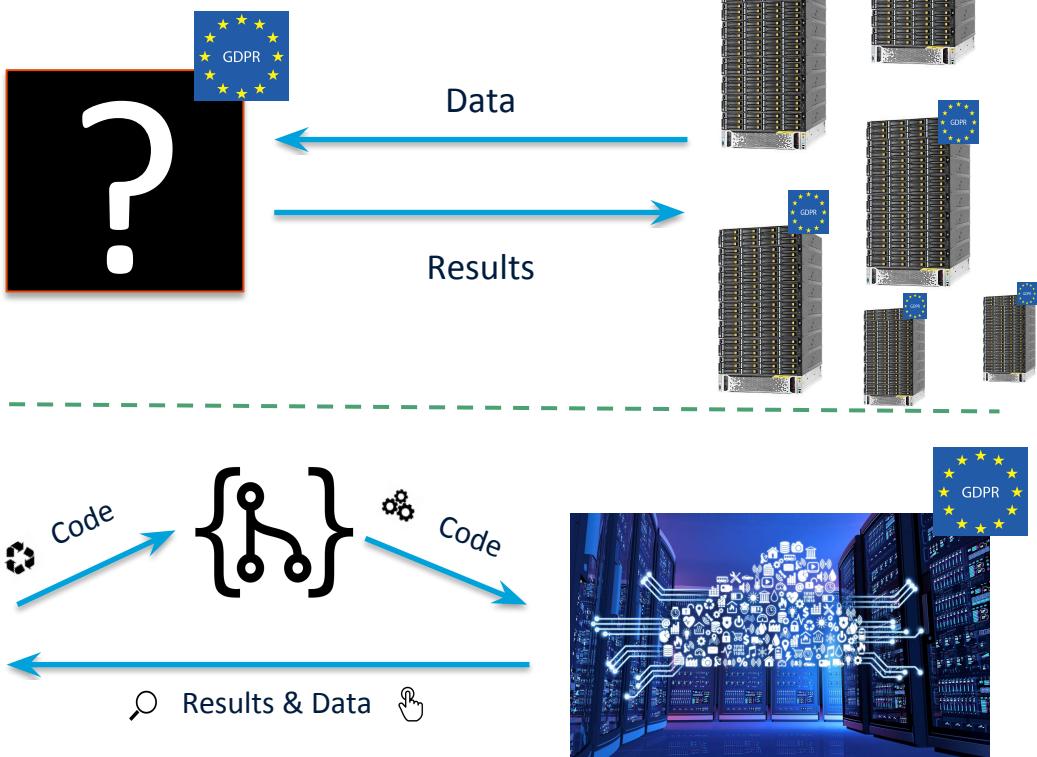
# Flexible computing

Current

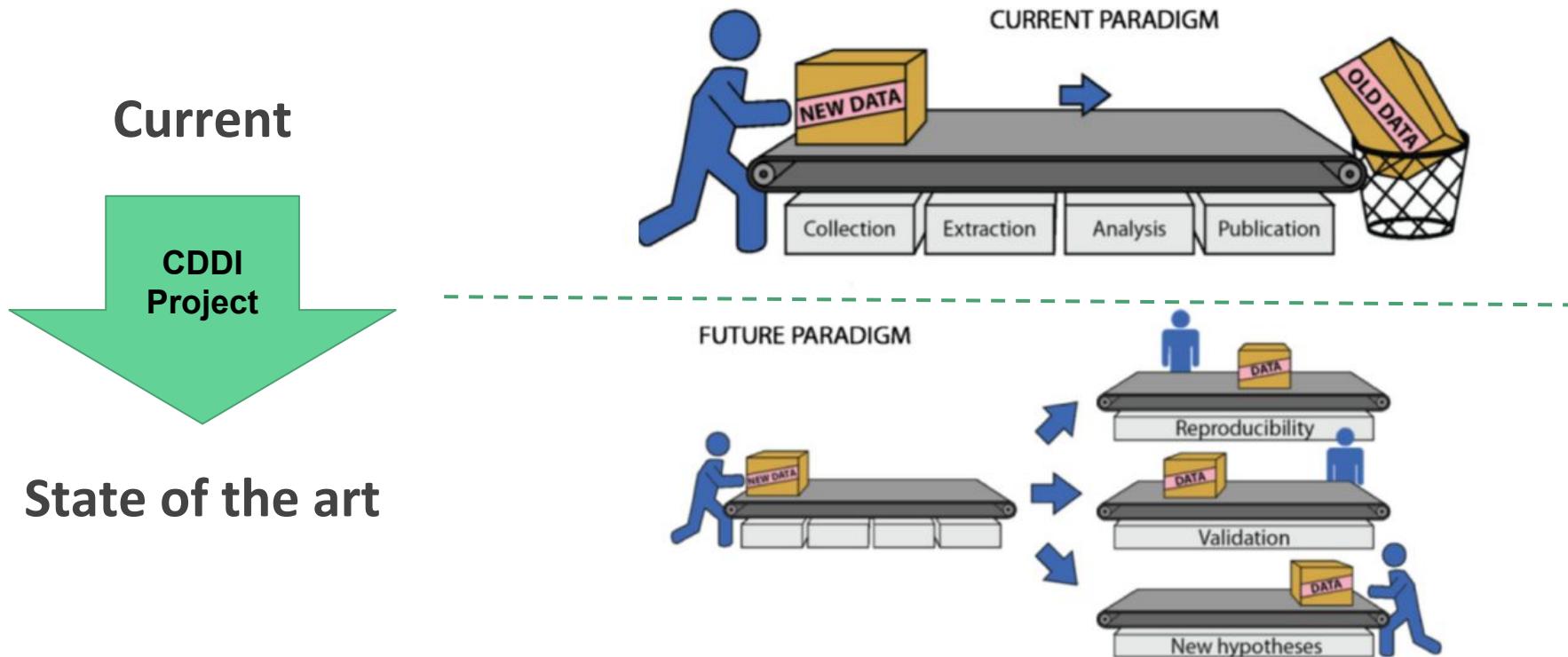


CDDI  
Project

State of the art

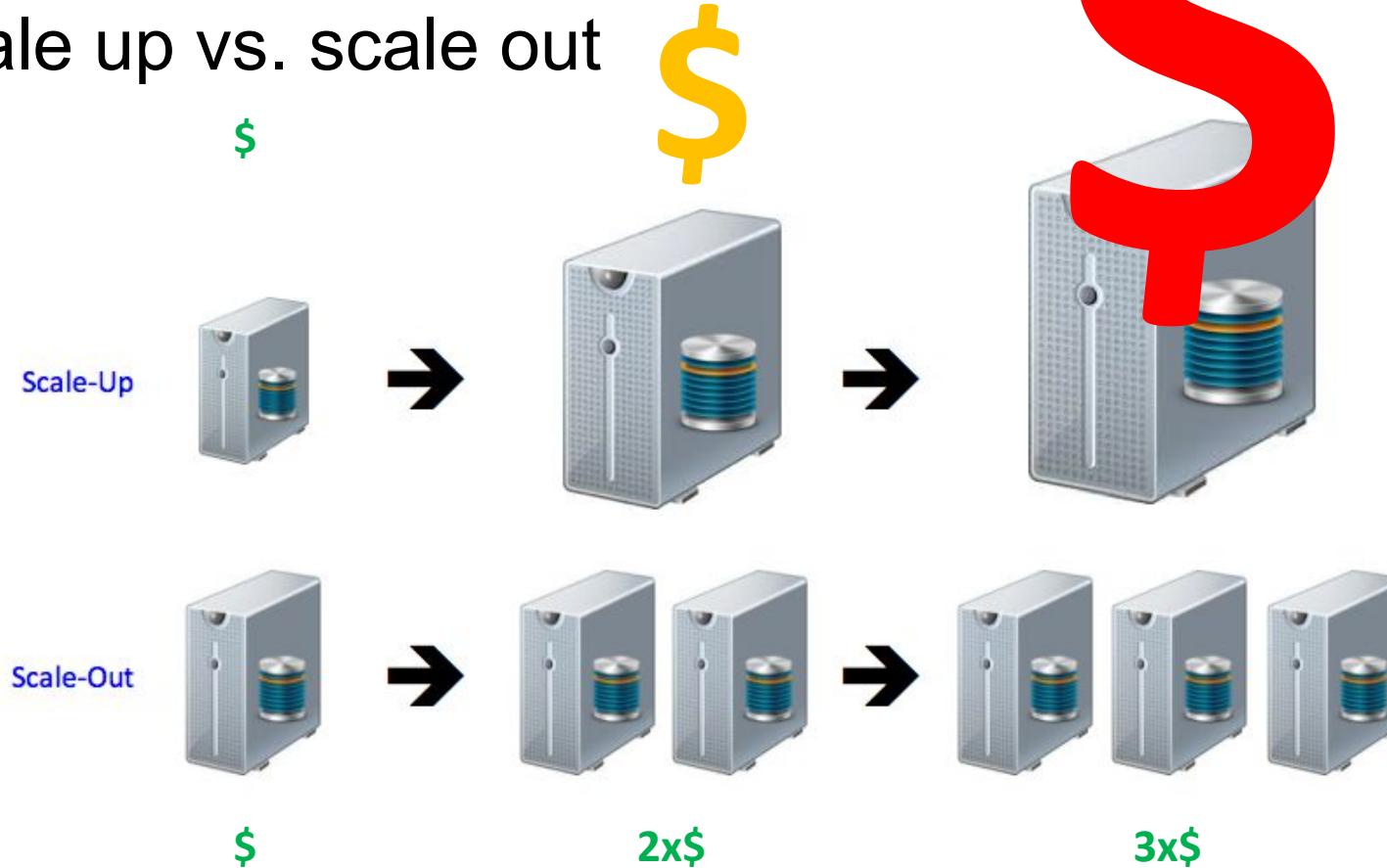


# Reproducible research



State of the art

# Scale up vs. scale out



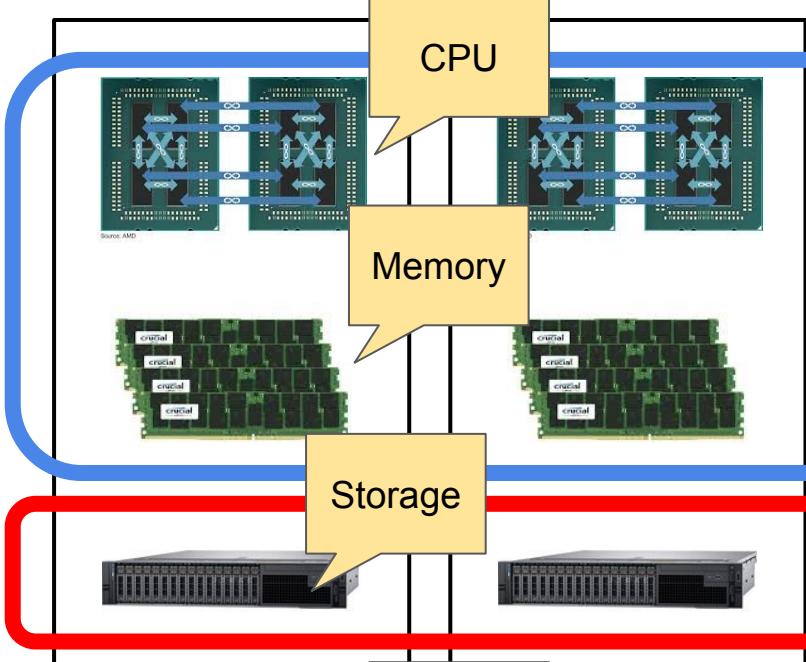
# What is DSRI?

- Infrastructure
  - shared pool of resources
  - high availability
  - collaborative
- Methods
  - reusable components (algorithms & code, data)
  - reusable data
  - reproducible methods
- Partnerships
  - Win-Win situation based
  - Dell, NVIDIA, AMD and MapR

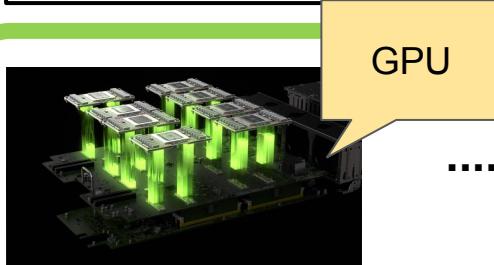
1

2

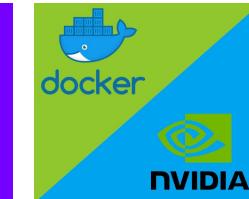
n

**kubernetes**

Big Data Hadoop &amp; Spark

**MAPR**

**RAPIDS**  
Open GPU Data Science

**kubernetes**

# DSRI project

3

## Performance indicators - numbers

	Base	Expanded
64 cpu cores per server	576 CPU cores	<b>960 CPU cores</b>
128 simultaneous threads per server	1152 threads	<b>1920 threads</b>
512 GB per server	4608 GB	<b>7680 GB</b>
120 TB net storage per server	1080 TB - 1350 TB*	<b>1800 TB - 2250 TB*</b>
DGX 1 - (8x Tesla V100 - 32GB GPU)		<b>1 x NVIDIA DGX 1</b>

\* estimated because of MAPR file system compression (assumed compression factor 2 - 2.5)

# MAPR & OKD (Origin Kubernetes Distribution)

Compute



OpenShift

Storage &  
Big Data



MAPR

# Why MAPR?

- Multi-Tenant
- High Available
- Scalable
- Secure
- Global Namespace
- Multi Location
- Multi Temperature
- Kubernetes integration
- NFS
- Open Source compliant

# Why OpenShift?

- Multi Tenant
- Open Source version of OpenShift
- Management interface that enables self service
- Kubernetes and Docker

# What can we do with it?

- IT view
  - Multiple servers running the same software
  - High available (20% of the cluster can be out of service, without interruptions or data loss)
  - Flexible (rolling upgrades, automated failover)
- Research view
  - Scalable (Hardware contributions through grants and investments)
  - Low maintenance (Managed through Research and IT)
  - Easy to use (Containers, Methods, Best practices and blueprints, support)
  - Easy to contribute
  - Easy to share and reproduce

# CDDI pilot project - 500 TB MRI data

## CDDI Pilot Project - Auditory Cortex

Alexander Malic

### Pilot project charter:

The goal of this project is to implement a Research data workflow which enables Researchers to share MRI images with other Researchers across and outside Maastricht University through a platform called XNAT.

### Scope:

#### In scope:

- Focus only on FPN MRI images for this pilot
- Set up a XNAT platform for sharing MRI data
- Develop a process which converts proprietary DICOM data images into open BIDS data

#### Out of scope:

- Fully functional XNAT workflow for all UM (idea is that extending the use case should be easily possible but might require changes and additional ethical approvals.

### Requirements:

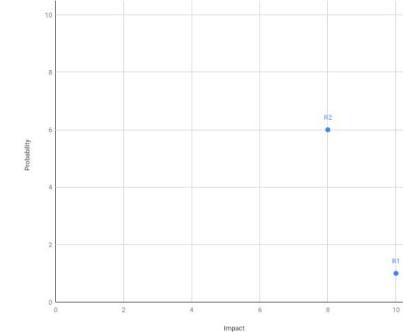
- DSRI Infrastructure in place
- 0.5 FTE for 6 months for Data engineer with MRI and XNAT expertise and ideally already familiar with BIDS data format (Jan 2019 to June 2019 → ~128 effective working days =  $0.5 * 128 * 8 = \sim 512\text{hr}$ )
- Approval and steering by CDDI

### Timelines:

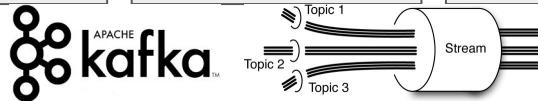
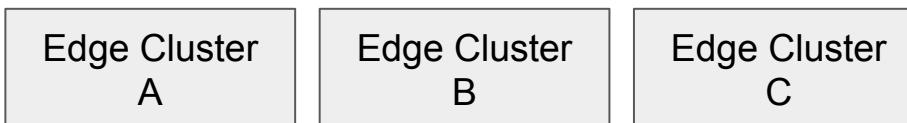
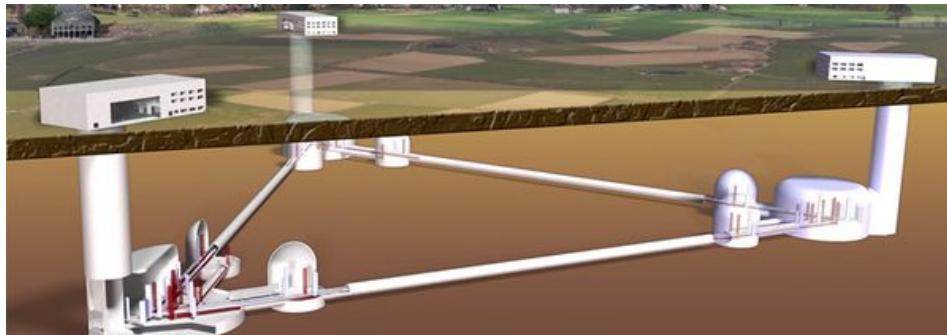
	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26
Pilot approved	♦																									
Project management																										
Resource planning																										
MRI Expert hiring																										
MRI Expert onboarding																										
Documentation																										
XNAT container and configuration																										
Dicom2Bids container																										
Bids2XNAT container																										
Pilot feasibility proved															♦											
Migrate and secure existing data																										
Pilot closure																										
Pilot project closed down																					♦					

### Risks

Risk	Impact	Probability	Description
R1	10	1	DSRI infrastructure not in place. Project will not start until infrastructure is available. Work on this pilot can start as soon as DSRI is online and does not need to wait for the DSRI project to be finished.
R2	8	6	MRI Expert with XNAT knowledge not available. If a MRI expert is not available for hiring, the project team will look into building up a colleague who is on PostDoc level. Project timelines need to be reevaluated in such a case.



# Einstein Telescope on DSRI as IoT device



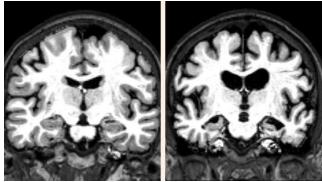
DSRI  
HA storage / scalable CPU & GPU



# BReIN project

## Maastricht - Backup Facility - FZ Jülich

### Randwyck



MRI data

```
@read1  
AGCTTATCCTCTGCTACCCCGGGTAGCGCACTTGATGTATTCAACAGC  
+  
BA1@C7CBCCC9C8; B2@>C?B@B@B=9?@B1 : AB7B?B8B?B6B. 7.  
@read2  
TTGGGCGGATCTCAGAACATGGATGTGATCCACACAGCATTCTG  
+  
?>?B@)<@, AA7A@C<?=@B; + )?B*@2=@=BB, -B6C>AB@B24  
@read3  
TATGCTAAGAAGGGGCTGATGAGTTGGTGTACGATATCACTGCCTC  
+  
A3AB: B1 : B; 9/0BBCBB<BB@AO?BB9: BB<@BB@7@6<@A@@@<3
```

Sequencing data

A	B	C	D	E
1	Use in Big Data analysis?	Study_ID	Location	Dataset
2	Yes	Study01	dixa data warehouse	DIXA-001
3	Yes	Study02	dixa data warehouse	DIXA-002
4	Yes	Study03	dixa data warehouse	DIXA-002
5	Yes	Study04	dixa data warehouse	DIXA-002
6	Yes	Study05	/ngs-data/data_storage/transcripcomics/ngs-data-storage_eastr	Carcinogenomics_eastr
7	Yes	Study06	http://ftp.biostorencdc.jp/w/DIXA-006	TG-GATES
8	Yes	Study07	/ngs-data/data_storage/transcripcomics/ngs-data-storage_eastr	STW_Magkoulopoulou
9	Yes	Study08	dixa data warehouse	DIXA-028
10	Yes	Study09	dixa data warehouse	DIXA-028
11	Yes	Study10	dixa data warehouse	DIXA-029
12	Yes	Study11	/ngs-data/data_storage/transcripcomics/ngs-data-storage_eastr	ESNAT
13	No	Study12	GSE31952	NTC_Van der Heijden
14	Yes	Study13	/ngs-data/data_storage/transcripcomics/ngs-data-storage_eastr	Deferme
15	Yes	Study14	E-MEXP-2458	Aunjaag_jennen
16	Yes		GSE33235	Deferme
17	Yes		GSE33235	Deferme
18				

Meta and other data

### Daalhof



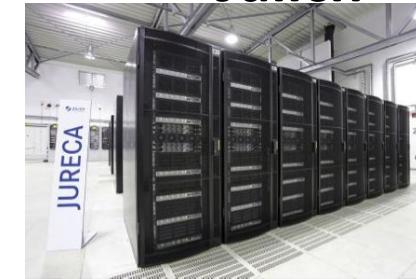
BReIN servers

### Heerlen



Backup facility

### Jülich



HPC

# BReIN - server specifications



Item	Description	Quantity
Compute and Storage node	Dell PowerEdge 7425 servers, 2x 32 core AMD Epyc (2x 64 core w. Rome?), 512 GB - 1TB Ram, 15x 12TB HDD	16
GPU node	1x NVidia DGX 1 (8x NVidia Tesla-V100 32GB), 512GB Ram	1



Item	Description	Quantity
Backup storage node	Baseline: 45-drives Storinator XL60 (max specs), 60x 12TB HDD	5

# 5. Vision



# Gitlab for DEVOPS → SCIOPS



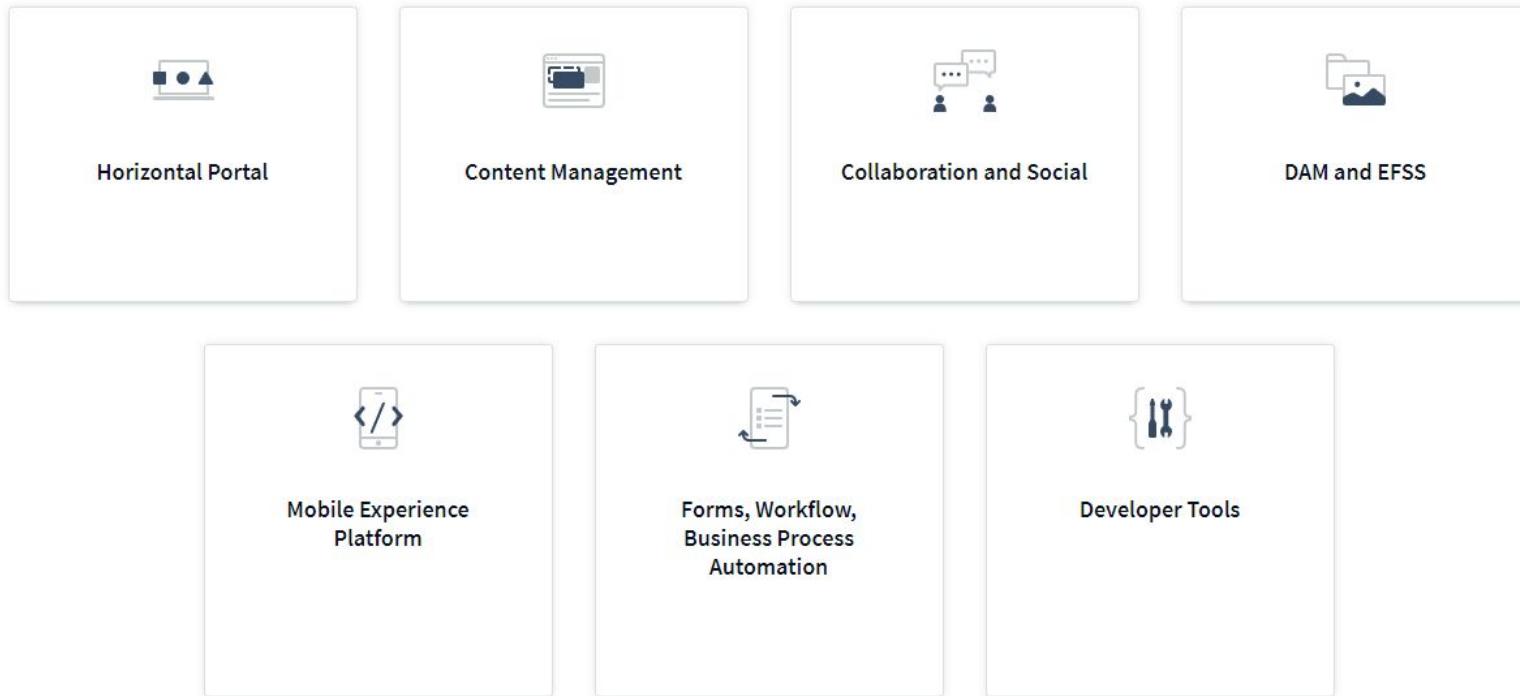
- Made for DEVOPS by DEVOPS
  - Case and Project management
  - Continuous Integration platform
- Container Registry/Cache
  - NPM & Maven Repository/Cache
  - One-Click deployment of code into containers



GitLab



# Liferay



# Open edX overview:

Built to showcase the latest in learning sciences and instructional design, the Open edX learner platform is driven by our community of developers, technology partners, research teams, and users.

## BUILT WITH

ubuntu  docker 

django  React 

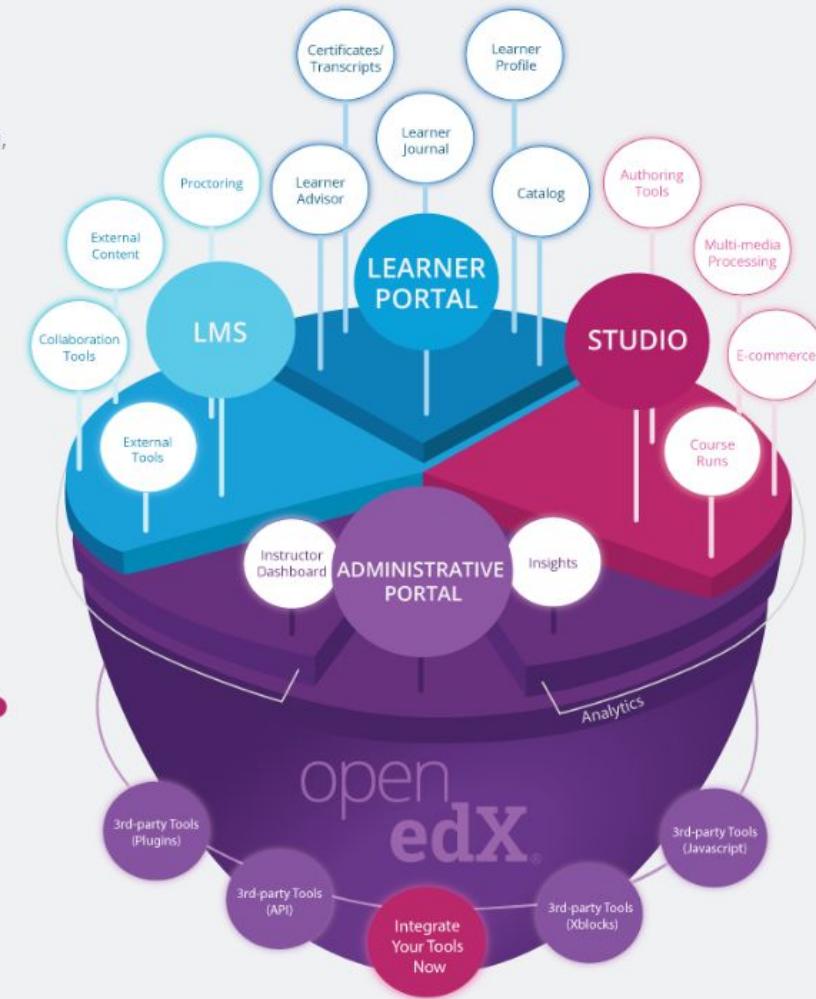
## RUNS ON

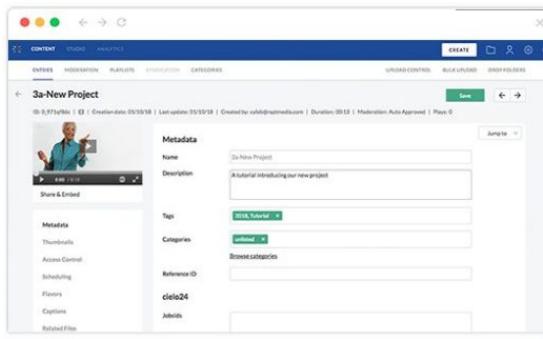
CLOUD 

Google Cloud  IBM Cloud 

Microsoft Azure  openstack. 

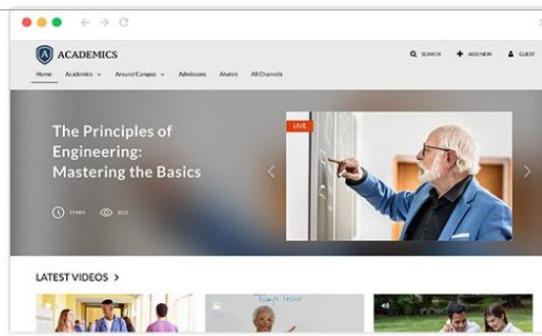
## ON-PREM





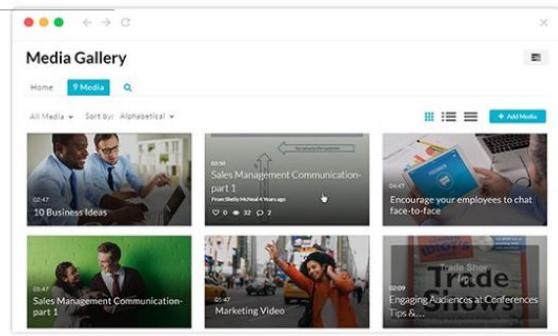
## Kaltura Video Platform

Manage and publish your media through one intuitive interface



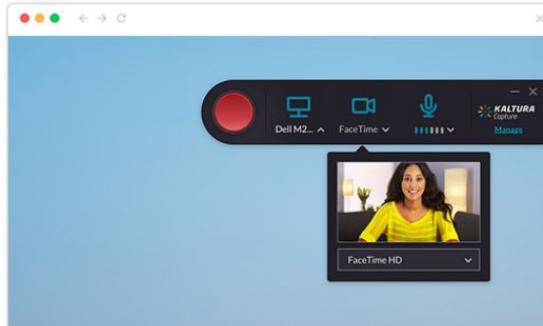
## Kaltura MediaSpace Video Portal

Create your own private YouTube-like video portal



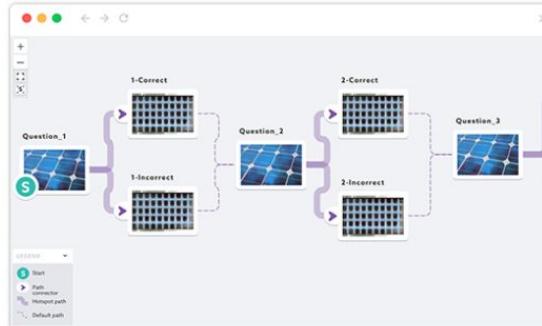
## LMS Video Plugins

Add a full layer of video to your learning management system



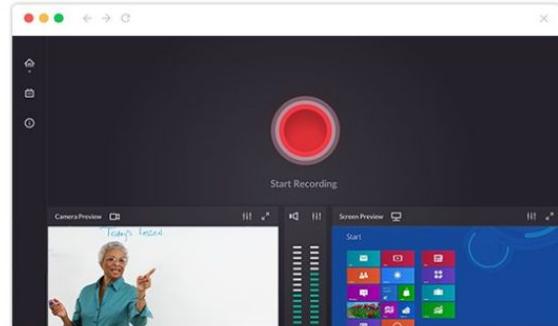
## Personal Capture

Give your end users powerful tools to easily create engaging videos from their own desktops



## Interactive Video

Create immersive video experiences for user engagement



## Lecture Capture

Manage all your lecture captures from any recording device



The background features a world map with a glowing blue aura around it. Numerous 'DSRI' labels are scattered across the map, appearing as bright blue lights or highlights on the landmasses.

# Thank You