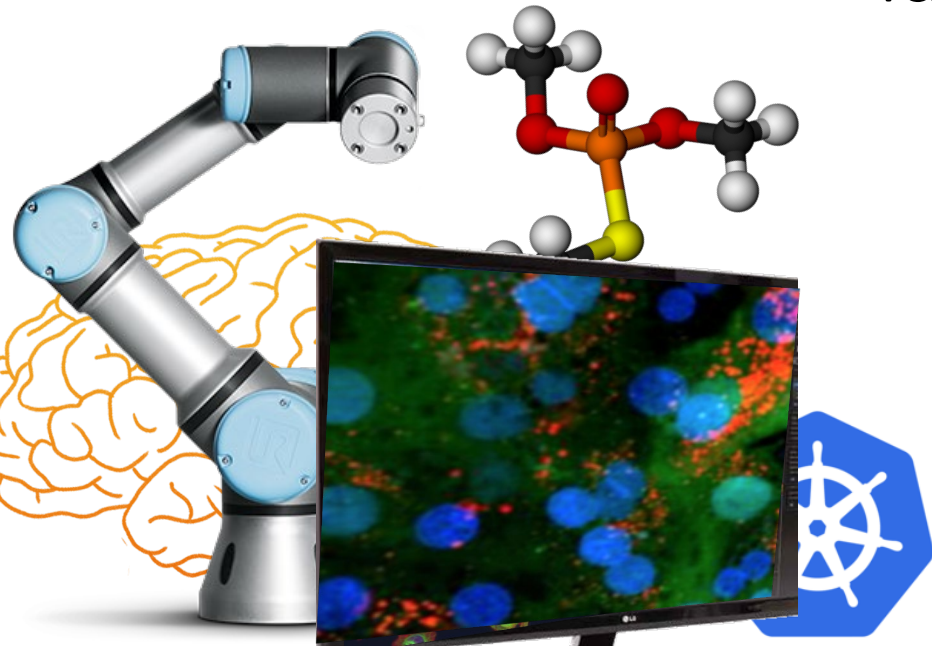


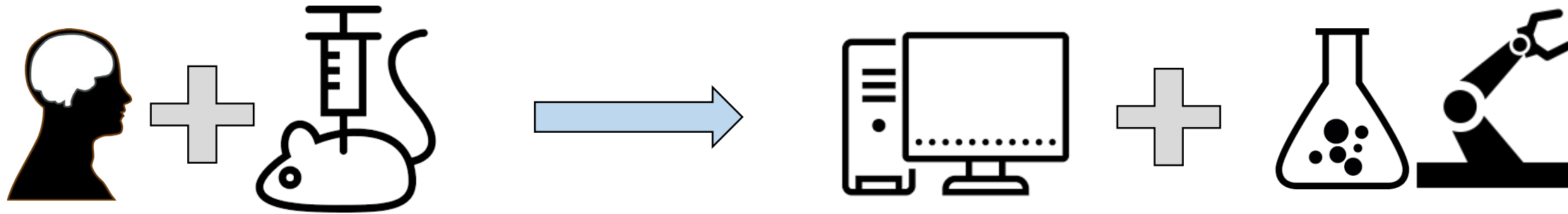
# Automating the process of continuously prioritising data, updating and deploying AI models in a robotised lab for drug discovery



Ola Spjuth <ola.spjuth@farmbio.uu.se>

Department of Pharmaceutical Biosciences, Uppsala University

[www.pharmb.io](http://www.pharmb.io)



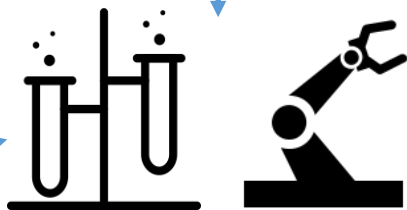
**Objective:** Accelerate drug discovery using AI and intelligent design of experiments.

- Predict safety concerns (fail early)
- Explain drug mechanisms
- Screen for new drugs

## Traditional hypothesis testing



Hypothesis



Experiments



Analysis and interpretation



Insight

- Iterative
- Flexible
- Mostly manual
- Slow

*revise*

*more*

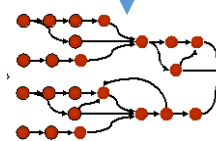
## Predictive modeling

- Retrospective analysis
- Hopefully predictive
- Expensive
- Limited for hypothesis testing

*Modeling and prediction*



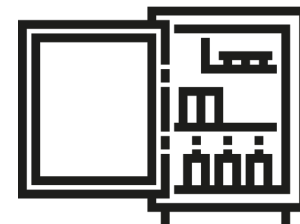
Prediction



Model

*query*  
*response*

*Data generation*

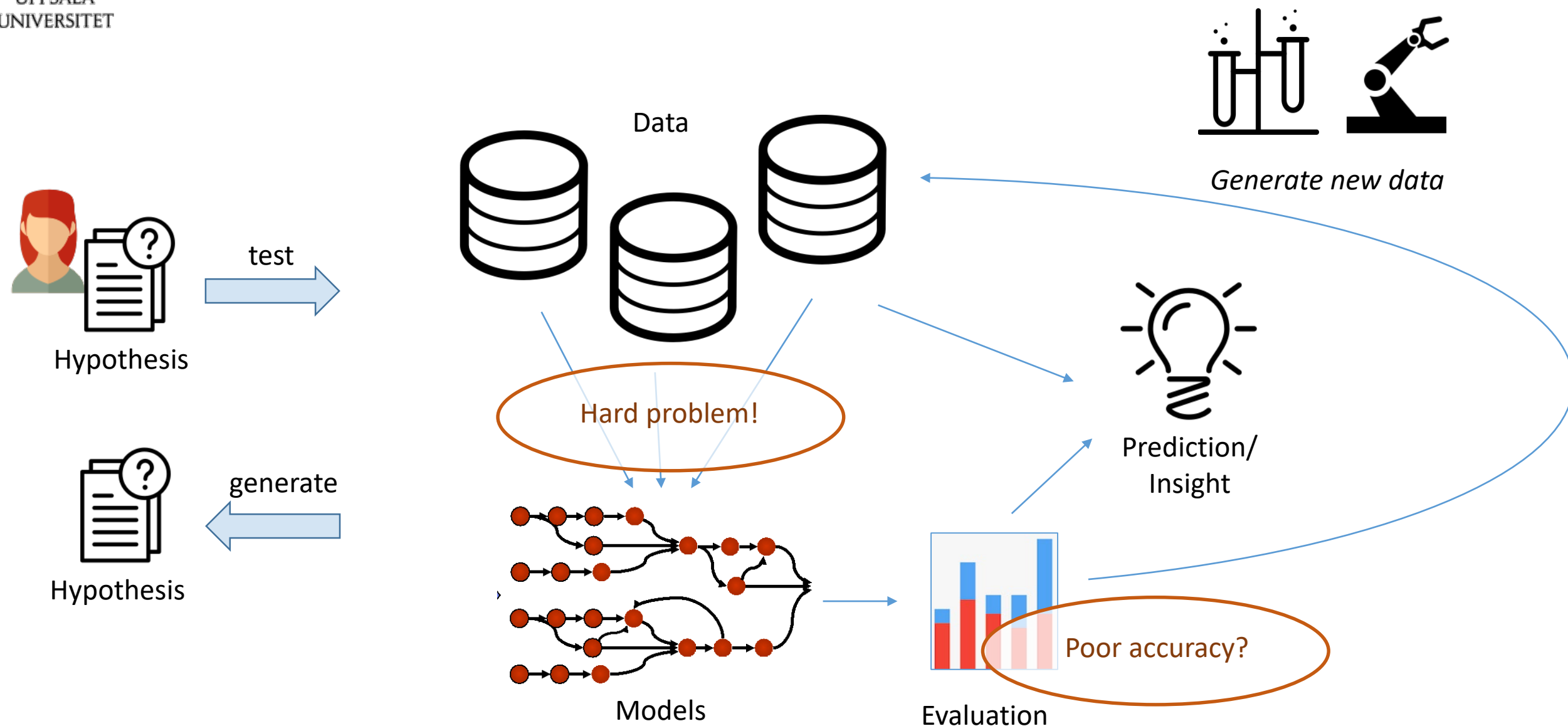


Database



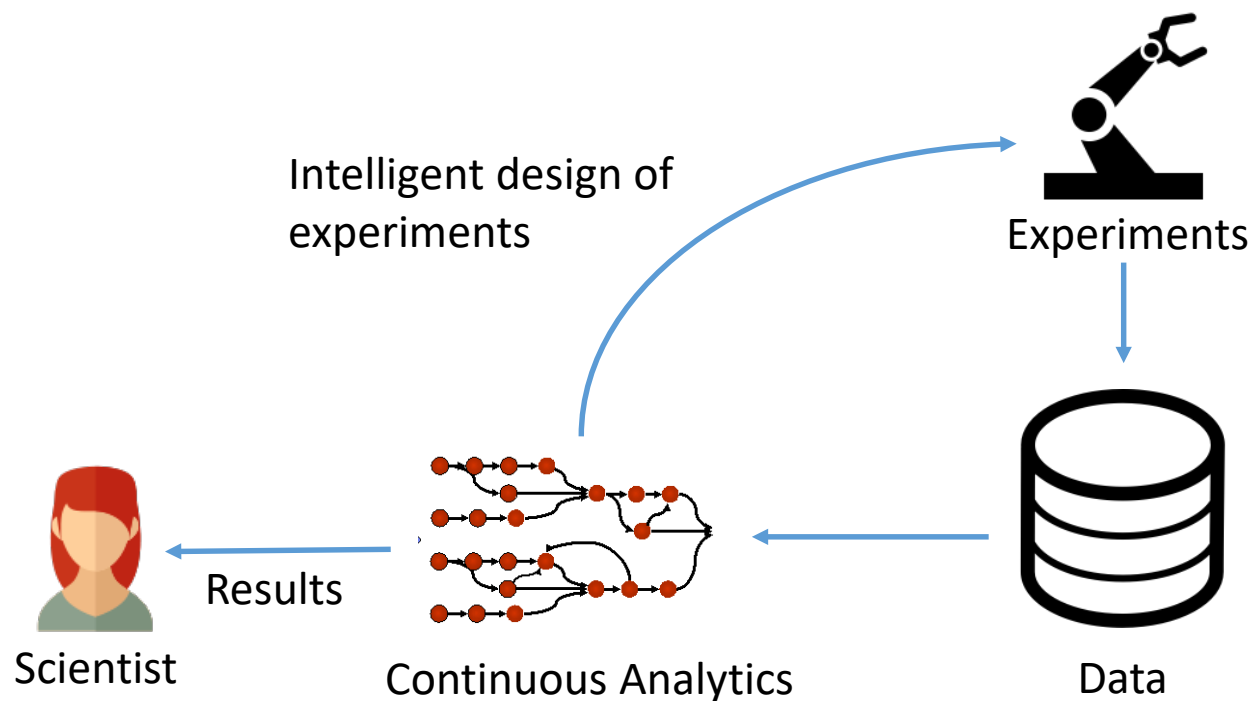
UPPSALA  
UNIVERSITET

# Data-driven science





# Intelligent experimentation



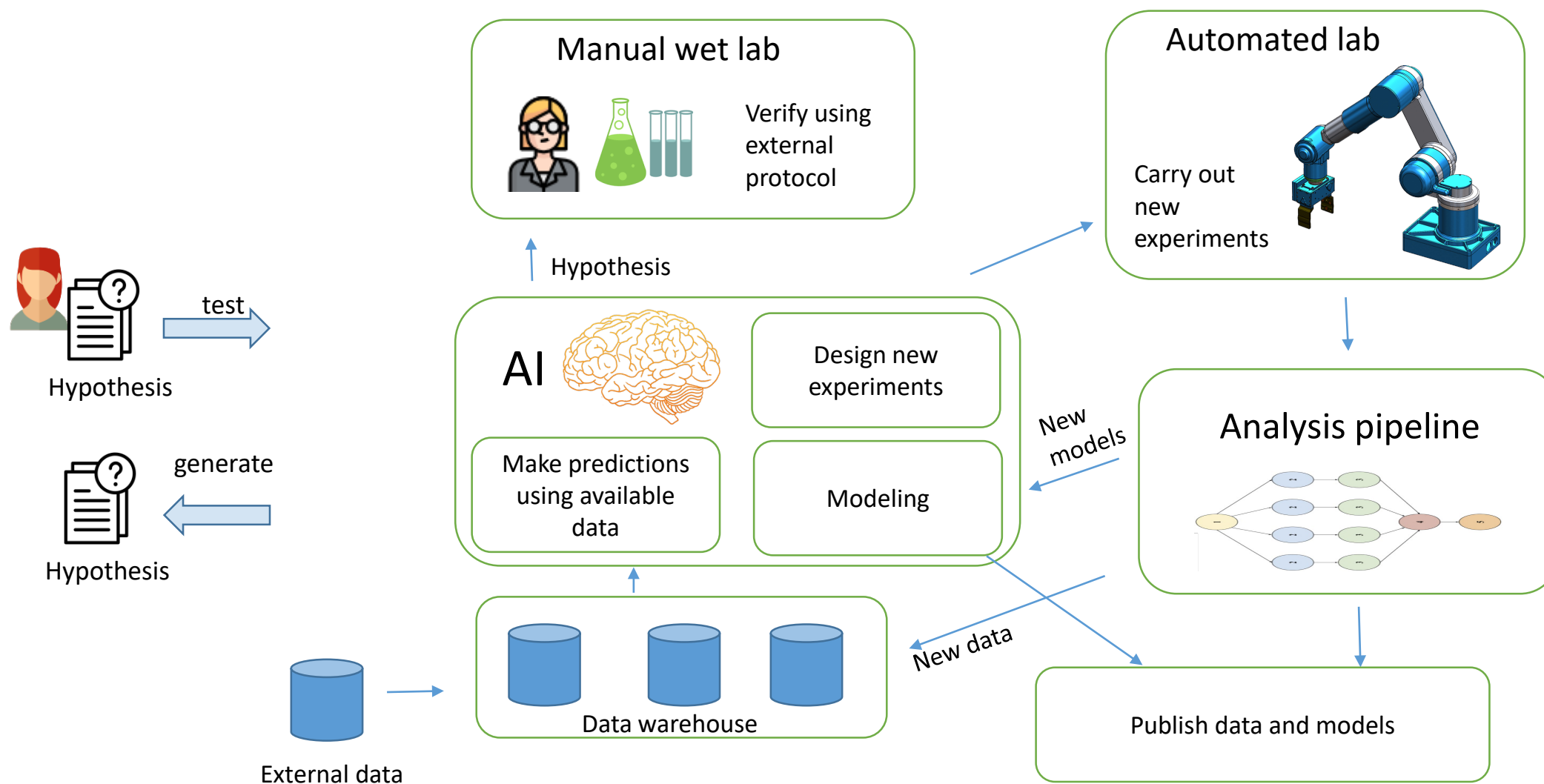
- What experiments should we do and how?
- Can we reduce search space?
- How store only interesting data?
- Can we replace experiments with predictions?

Informatics system



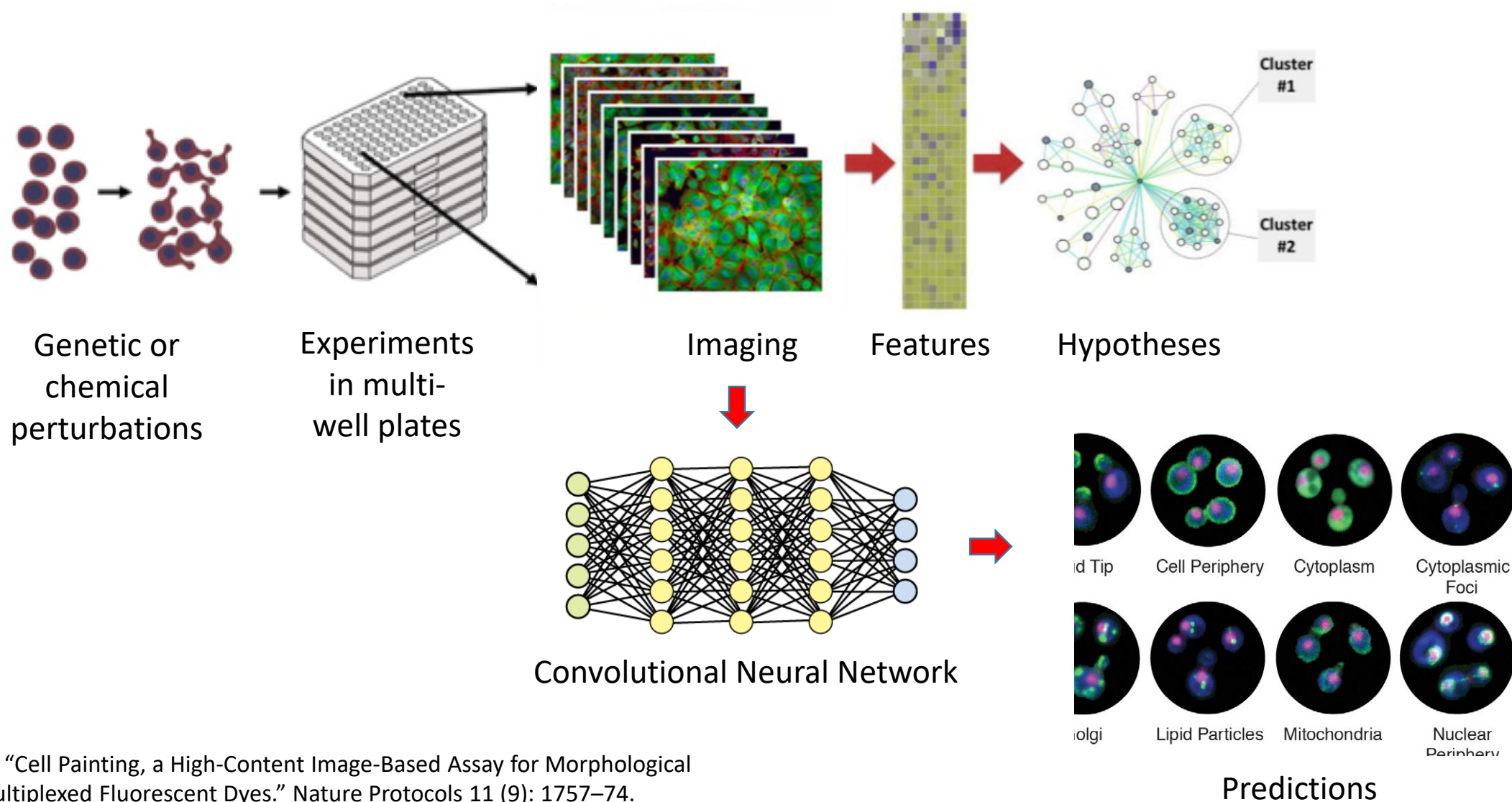
UPPSALA  
UNIVERSITET

# Vision: Intelligent systems for assessing drug leads





# High-content cell profiling

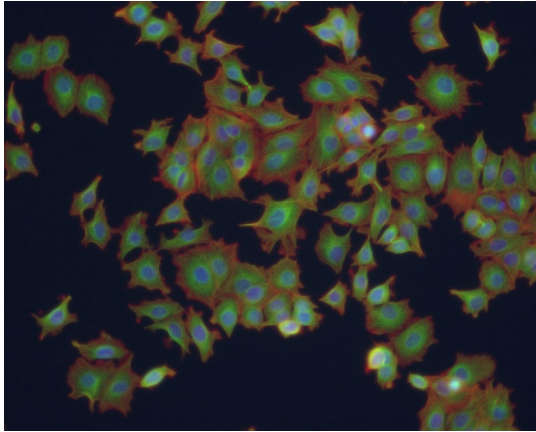




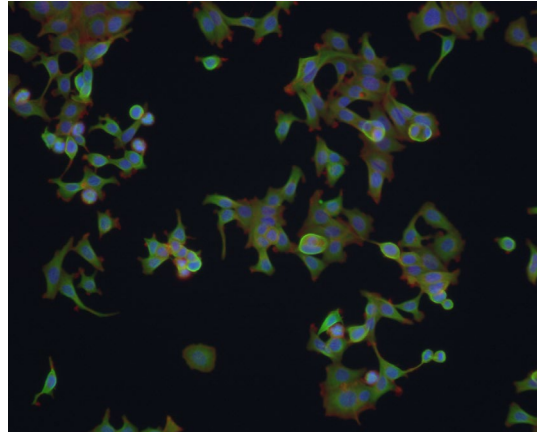
UPPSALA  
UNIVERSITET

# Classify images into biological mechanisms

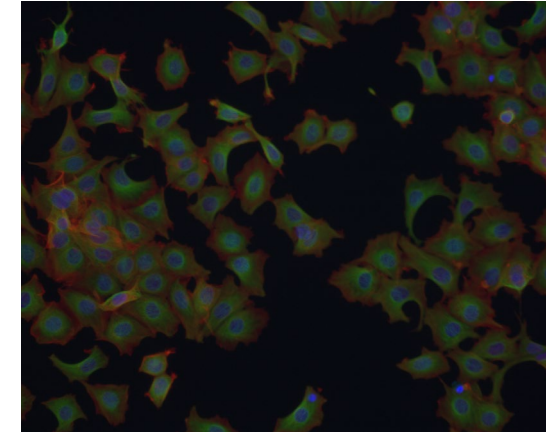
Protein degradation



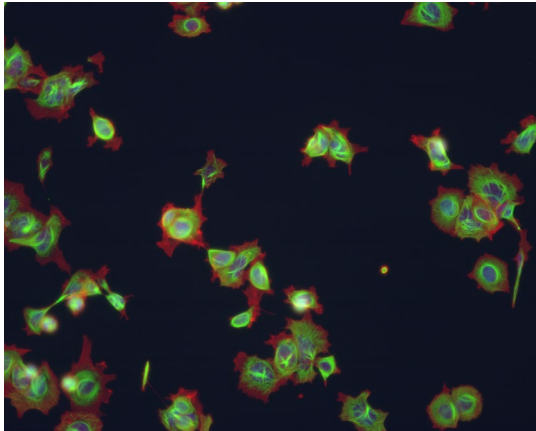
Cholesterol-lowering



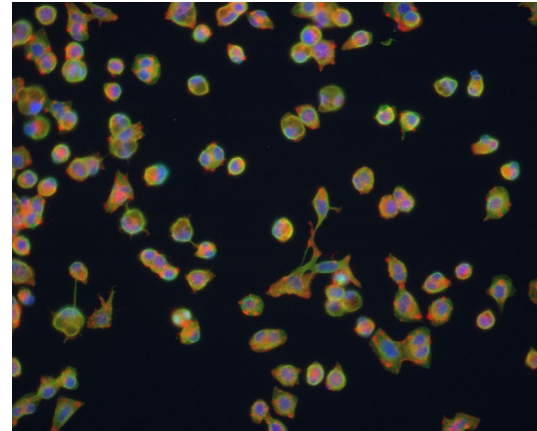
DNA replication



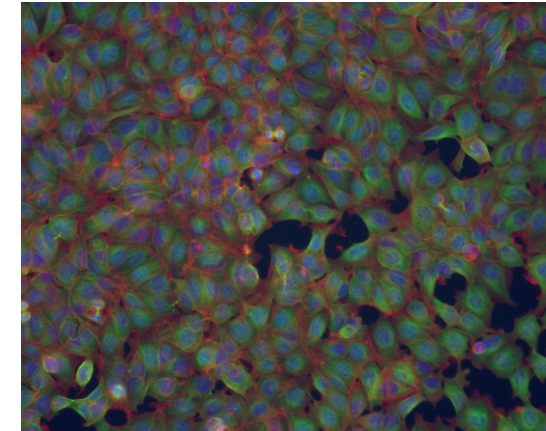
Microtubule stabilizer



Actin disruptor



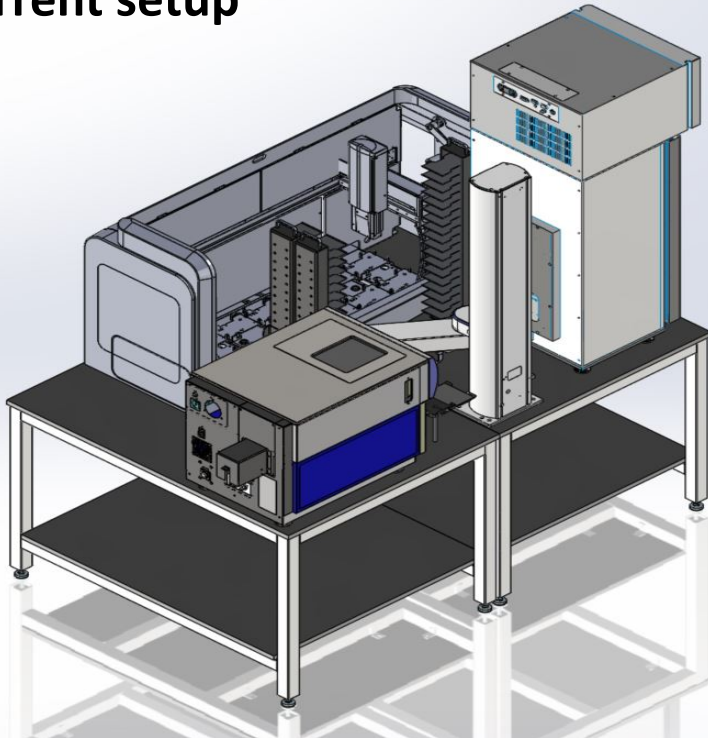
Kinase inhibitor



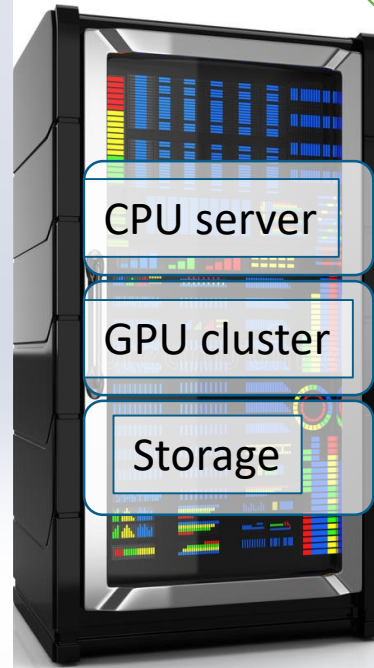


# Automated cell-based experiments

## Current setup



Automated plate handling



Informatics system

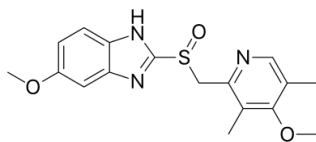
## Open source robotized cell lab



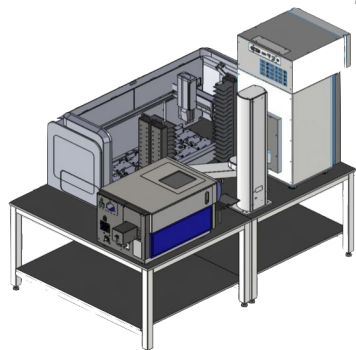
Automate more protocols



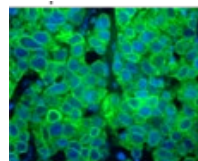
UPPSALA  
UNIVERSITET



chemicals



Robotized lab



images



cell profiles

Chemical DB



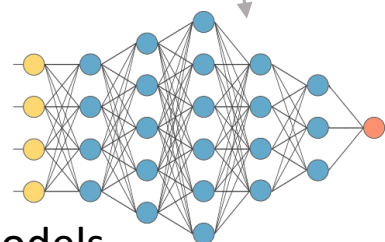
link



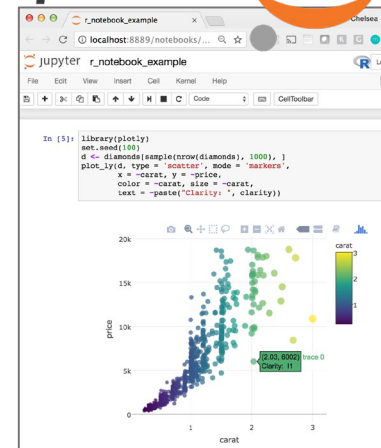
OMERO



CPU/GPU/HPC cloud



Models



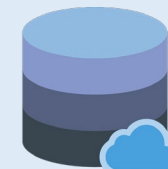
Notebooks



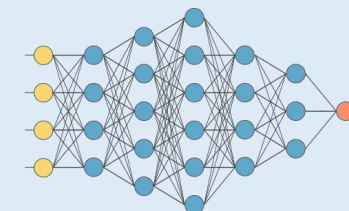
External  
user



Services



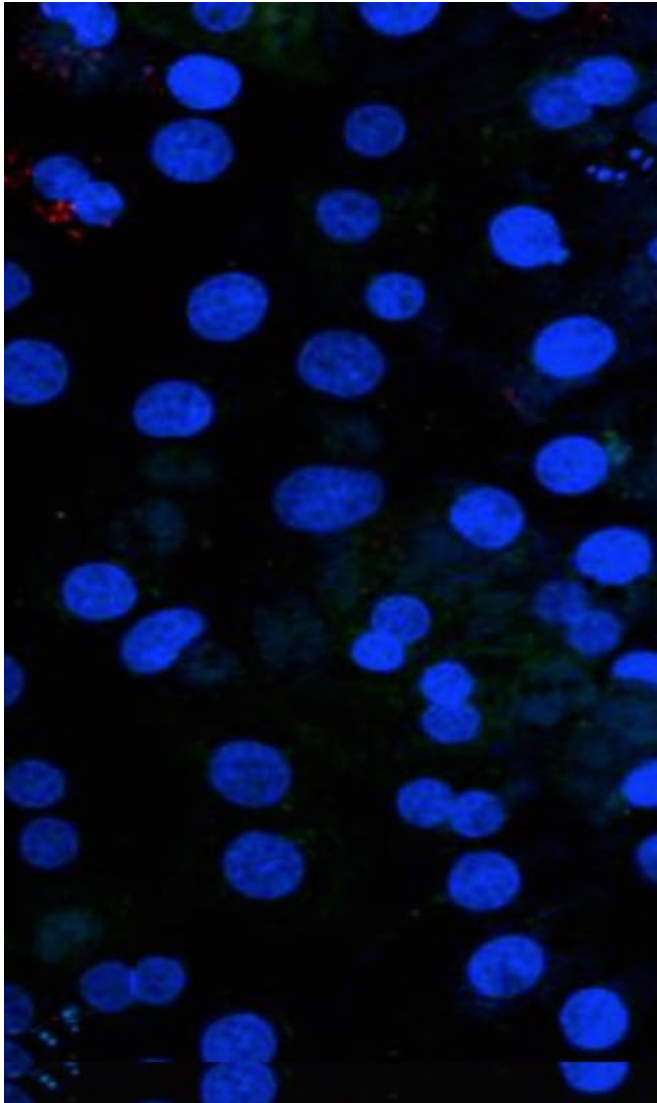
Data



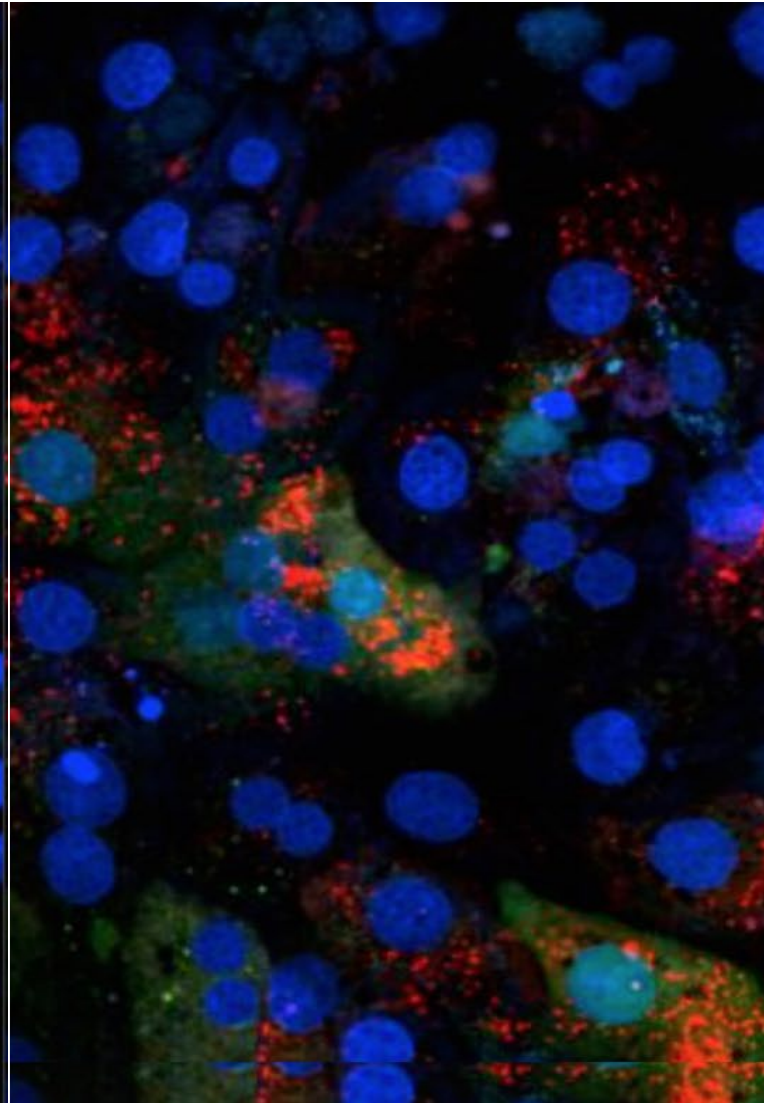
Models

Public services

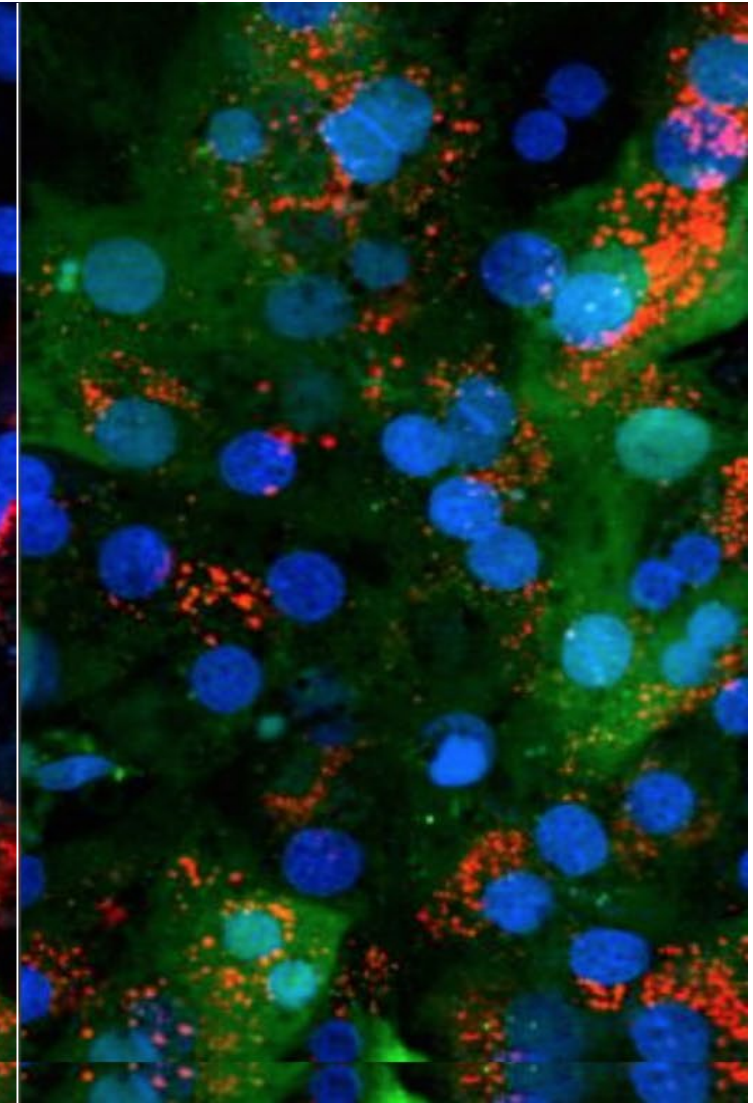
- Fluorescent LNPs (lipids)
- Fluorescent Cargo (mRNA)
- Fluorescent Product (protein)



No  
LNPs



Partial LNP  
uptake

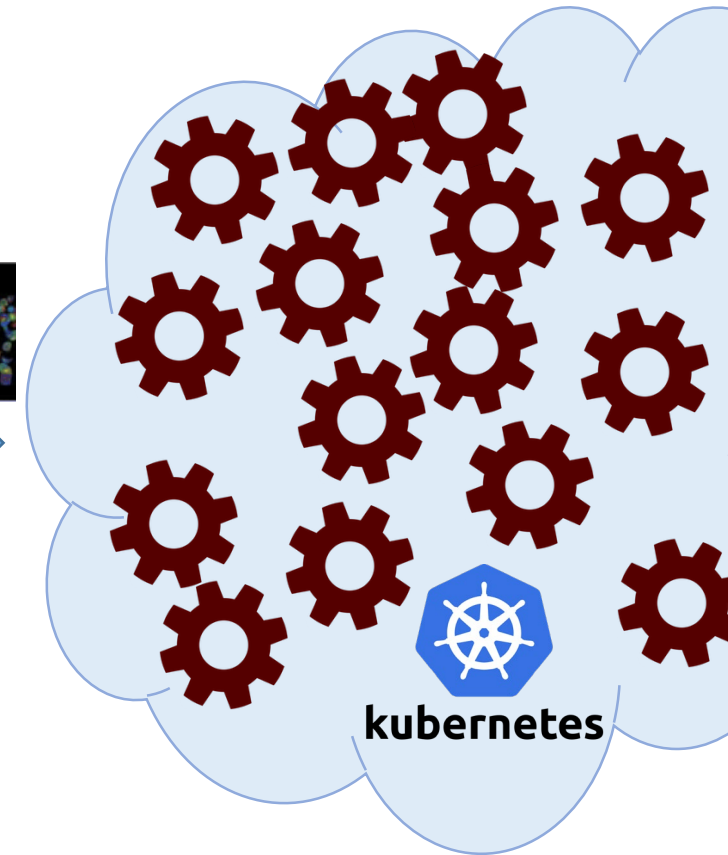
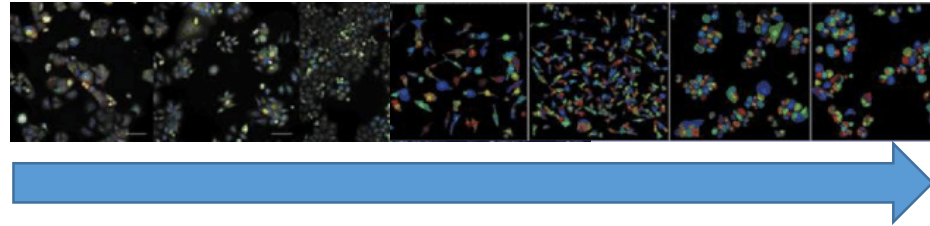
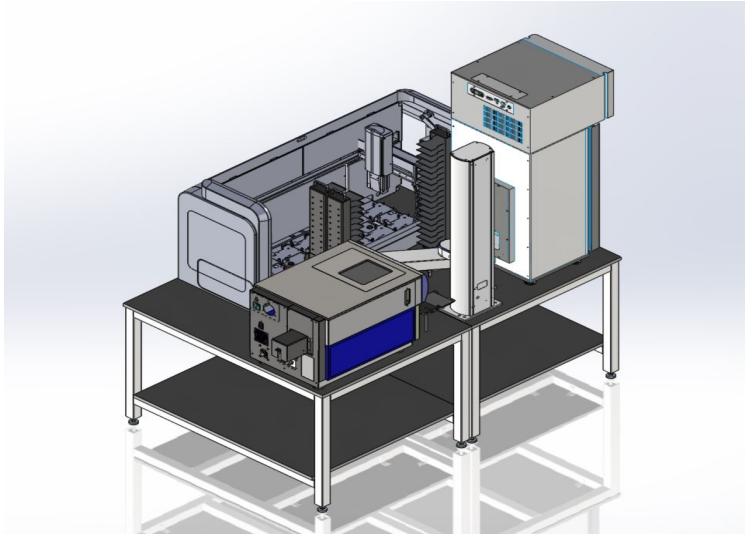


LNP uptake and mRNA  
decoding





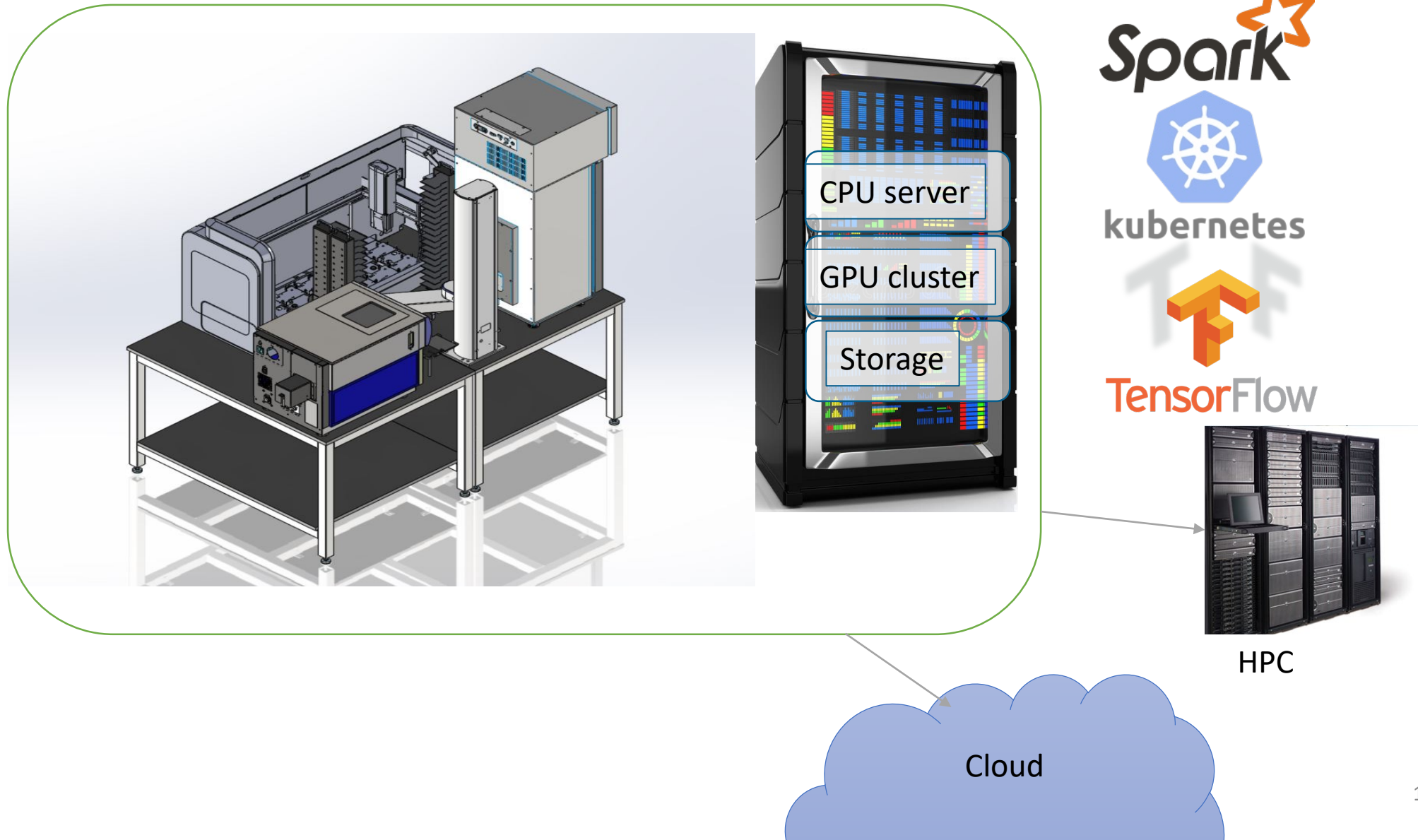
# Dealing with large scale data



- High volume, high velocity
- Continuously process data, train models, serve models
- Embrace scalable virtual infrastructures (cloud) and microservices (containers)
- Intelligently prioritize what to store



# Designing a flexible and scalable informatics system

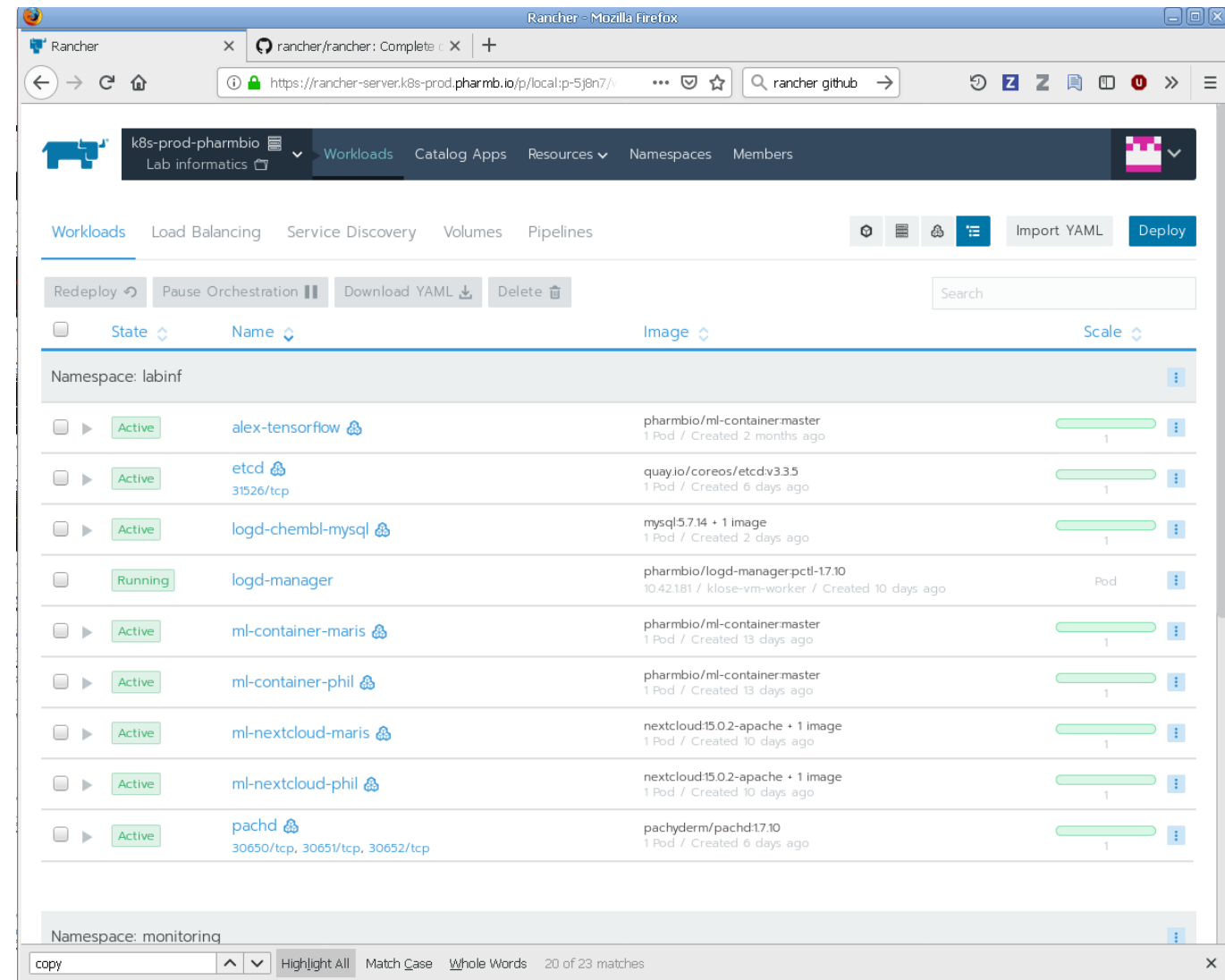




# Multi-cluster Kubernetes management

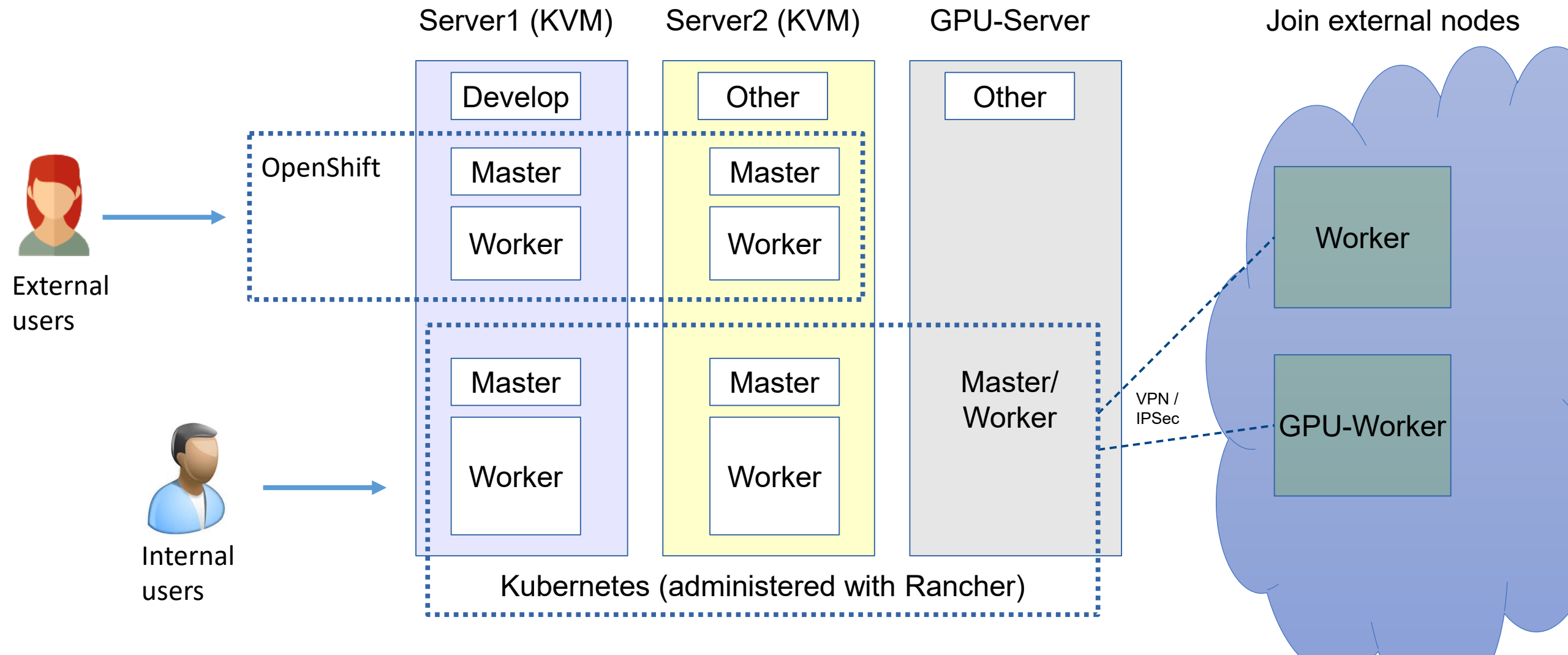


- Kubernetes Setup and admin in the cluster
- High availability. Kubernetes extras such as Networking, Helm, Nginx-load-balancer with tested and matched versions.



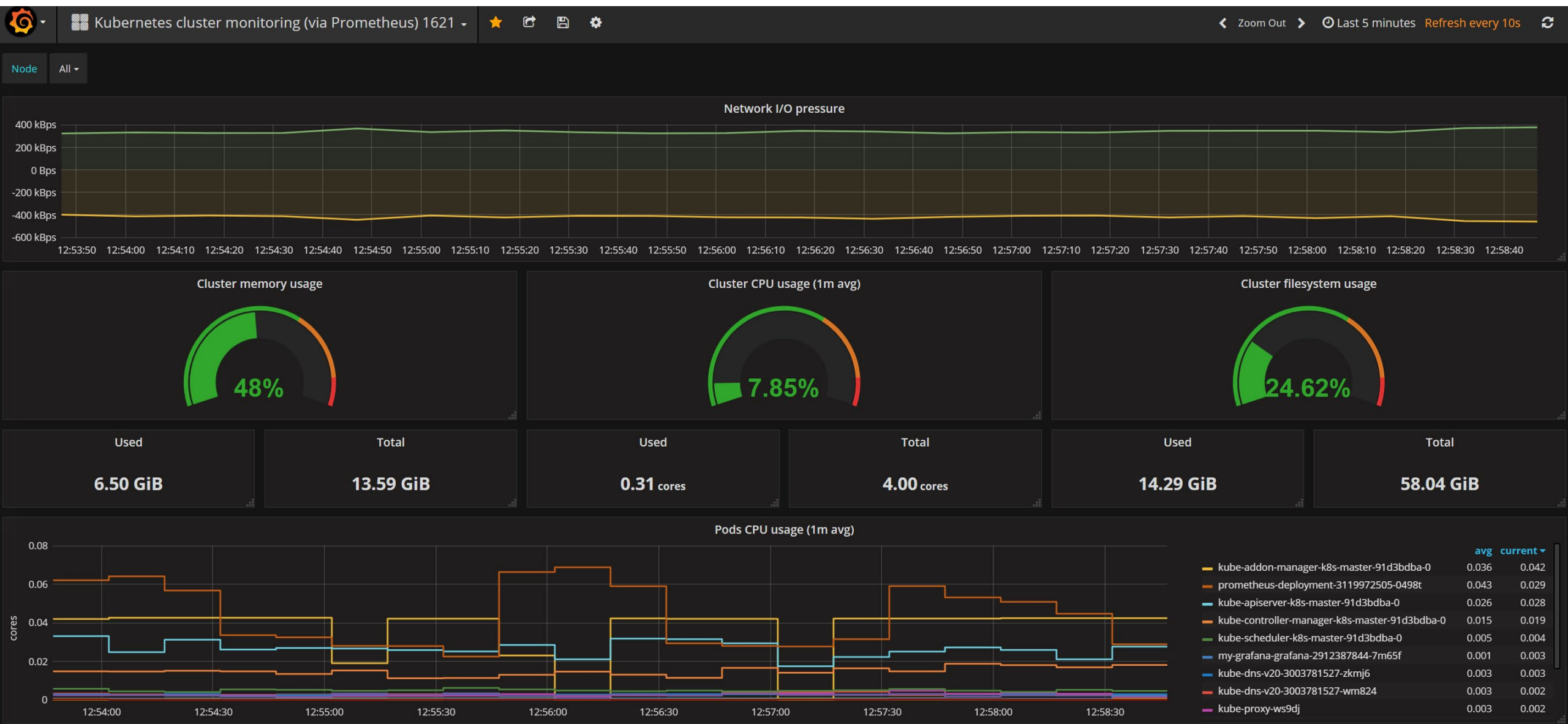


# Hybrid infrastructure



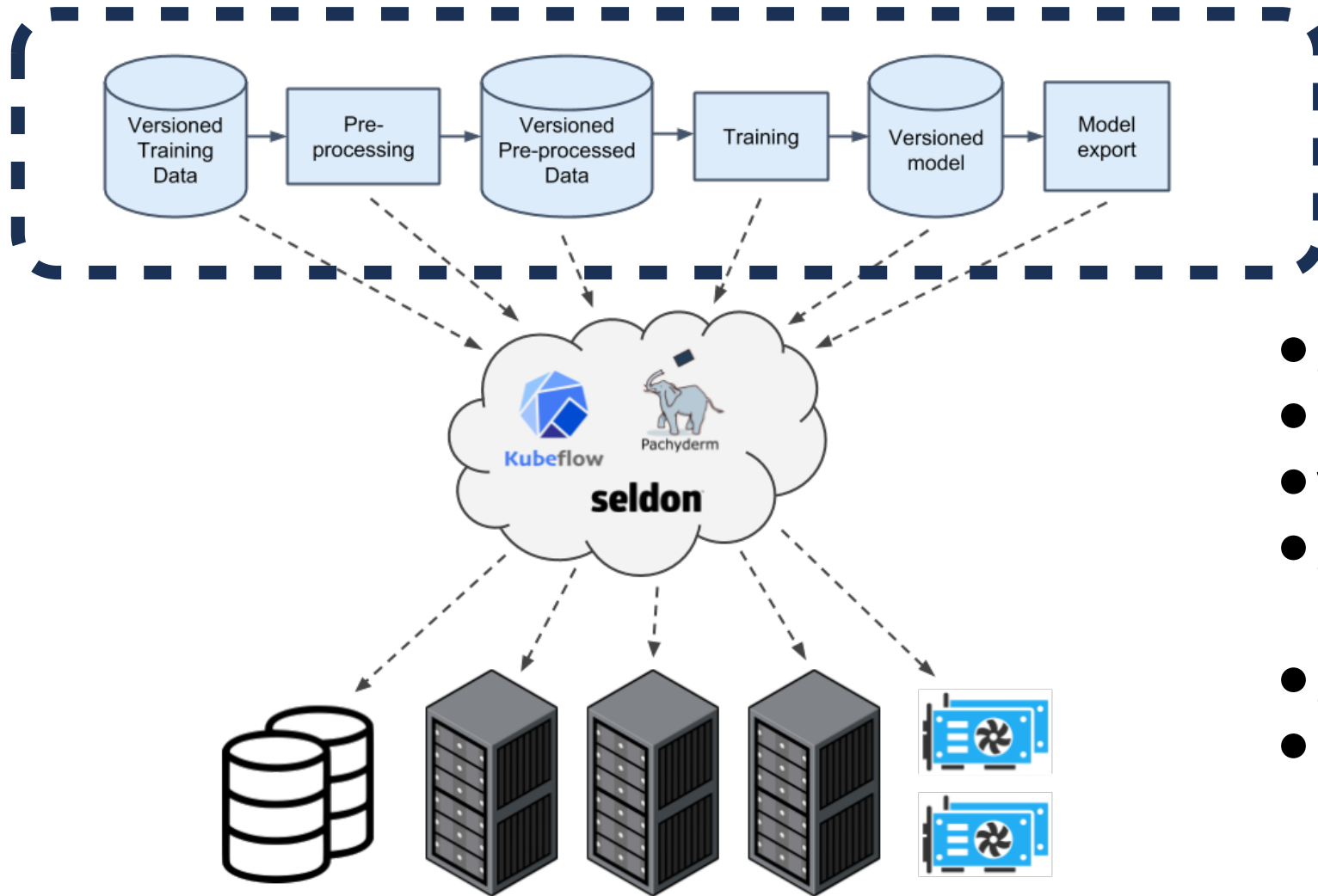


# Logging and monitoring





# AI/Machine Learning life cycle



- Structured processing steps
- Governance
- Versioning
- Streamline analysis on high-performance e-infrastructures
- Support reproducible data analysis
- Enable large-scale data analysis

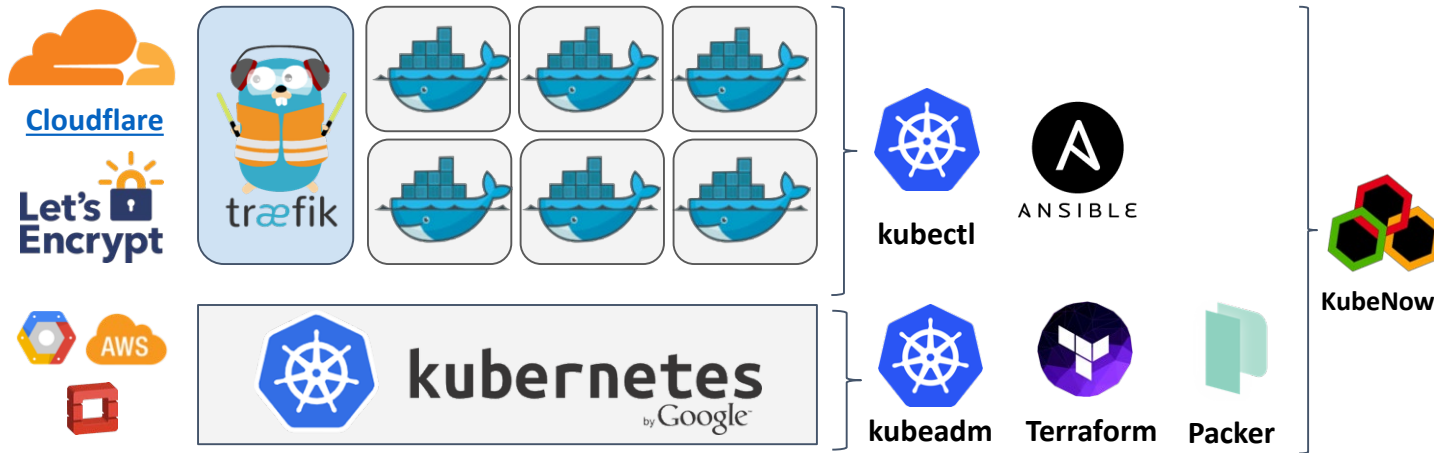


# Tools we develop and use

<https://github.com/kubenow>



- Easy deployment of virtual infrastructures on IaaS
- Containerize tools, orchestrate microservices with workflow systems on top of Kubernetes and OpenShift



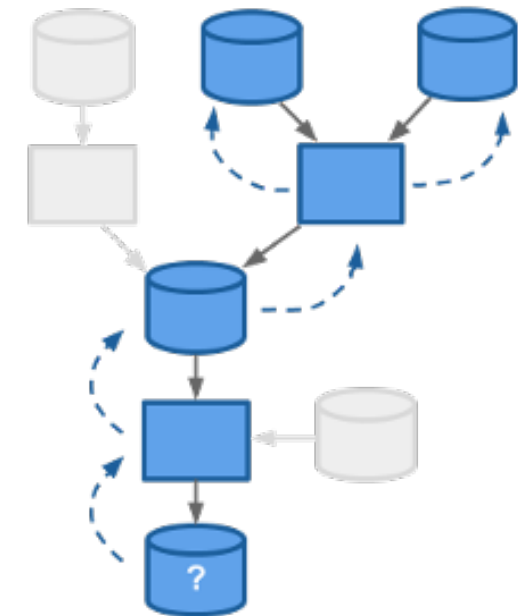
Capuccini et al. **On-Demand Virtual Research Environments using Microservices.** <https://arxiv.org/abs/1805.06180>

<https://www.pachyderm.io>



Pachyderm

- Data pipelining and data versioning layer for Kubernetes



Novella et al. **Container-based bioinformatics with Pachyderm.** *Bioinformatics.* 35, 5, 839-846. (2018). DOI: <http://dx.doi.org/10.1093/bioinformatics/bty699>



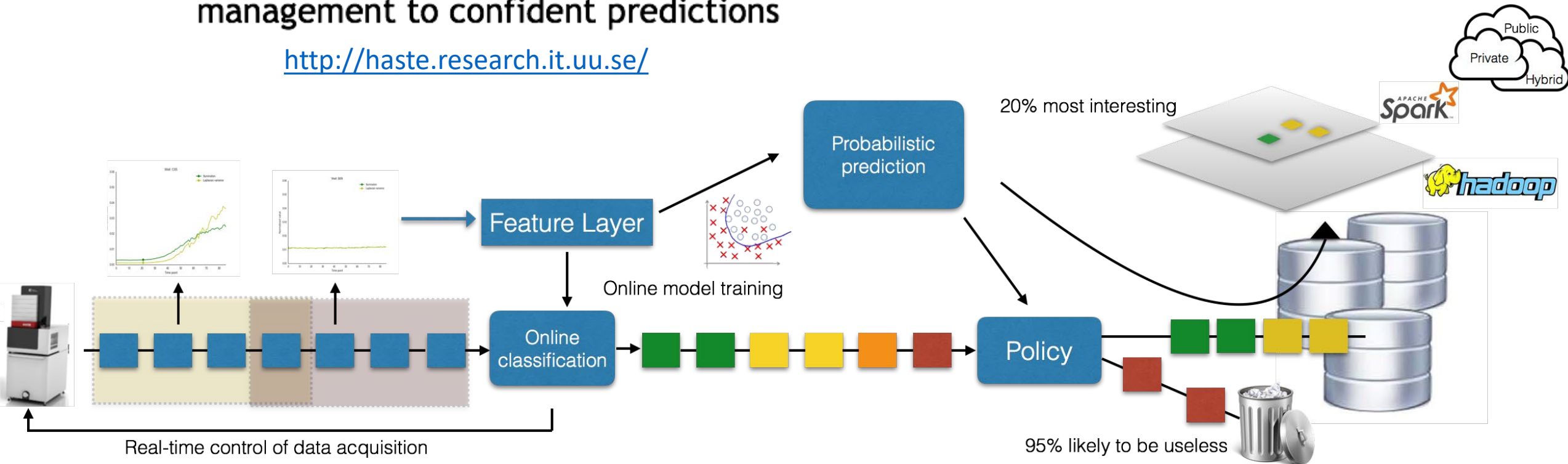
STIFTELSEN för STRATEGISK FORSKNING



# HASTE: Hierarchical Analysis of Spatial and Temporal data

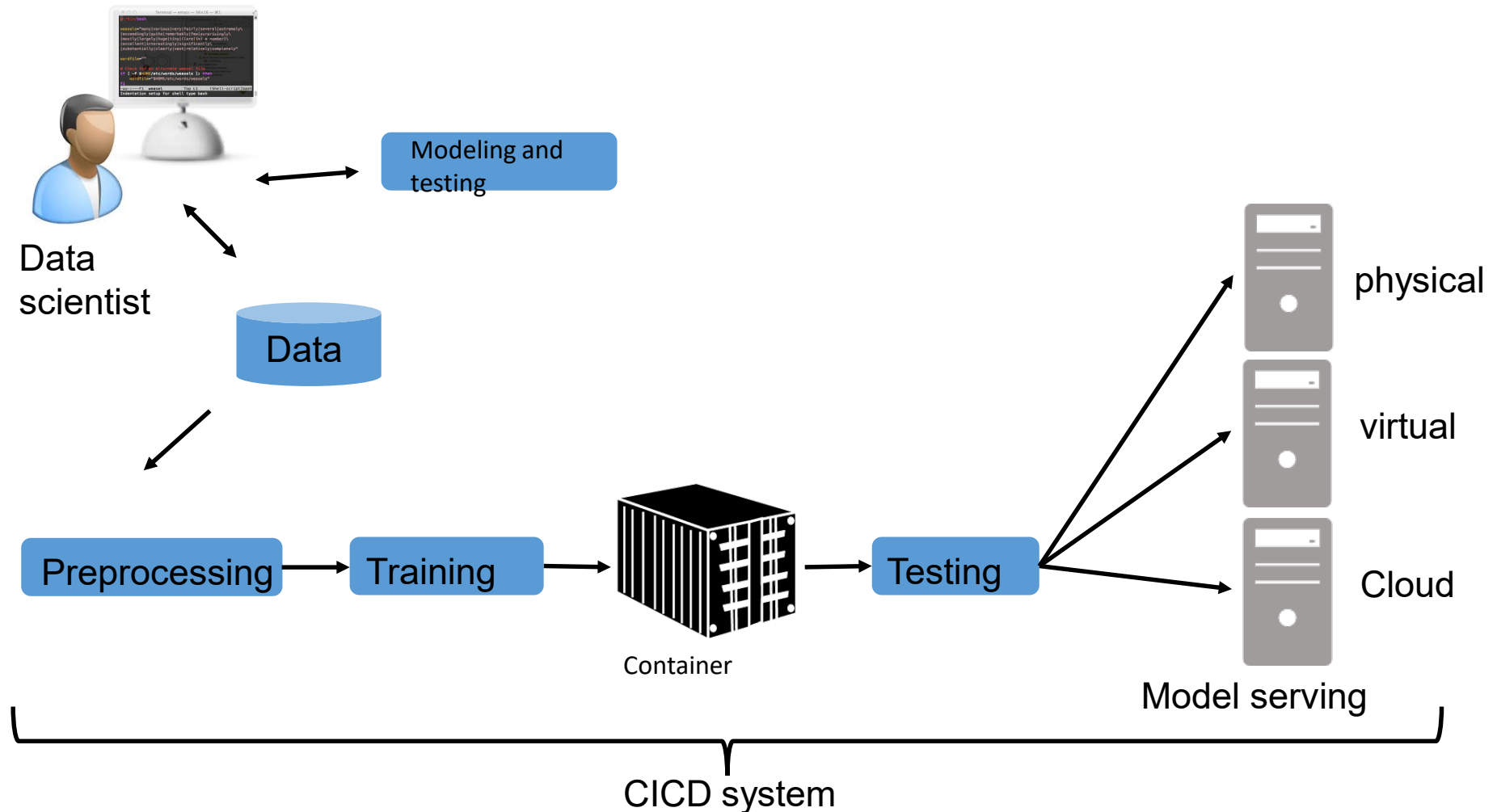
- from intelligent data acquisition via smart data-management to confident predictions

<http://haste.research.it.uu.se/>



# Continuous analytics

- Empower data scientists and automate from data to deployed models





# FAIR data and services

## FAIR<sup>1</sup>

- **Findable**
  - Semantic service discovery (JSON-LD)
- **Accessible**
  - Web UI
  - Programmatic API (OpenAPI)
- **Interoperable**
  - Open API
  - Semantic annotations (JSON-LD)
- **Reproducible**
  - Published scientific workflow (Pachyderm)

## Services

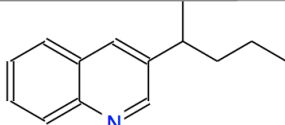
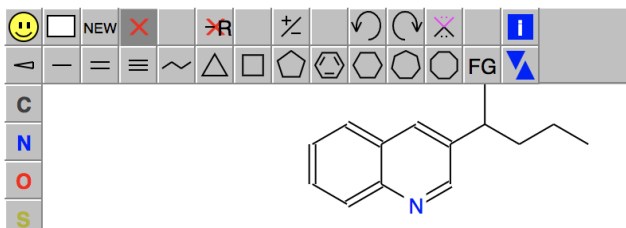
- Portable (Containers)
- Resilient (Kubernetes)
- Scalable (Kubernetes and cloud computing)

<sup>1</sup> Wilkinson, Mark D., et al. "The FAIR Guiding Principles for scientific data management and stewardship." *Scientific data* 3 (2016). <sup>22</sup>



# Examples of what we serve

## Target (safety) profiles

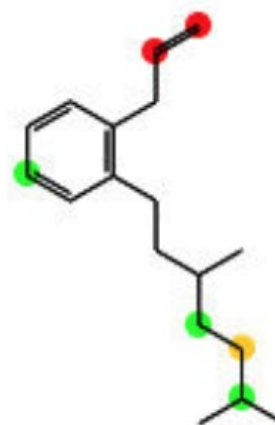


ACHE	ADORA2A	ADRB1	ADRB2	AR	AVPR1A	CCKAR	CHRM1	CHRM2
CHRM3	CNR1	CNR2	DRD1	DRD2	EDNRA	HTR1A	HTR2A	KCNH2
LCK	MAOA	NR3C1	OPRD1	OPRK1	OPRM1	PDE3A	PTGS1	PTGS2
SCN5A	SLC6A2	SLC6A3	SLC6A4					



**ADORA2A**  
p(N)=0.939

## Site-of-metabolism and reaction types



1. Input a drug candidate

2. List all potential reactions after probability of reaction to occur:

88.5% Reaction type 1  
75.2% Reaction type 2  
:  
0.0% Reaction type N

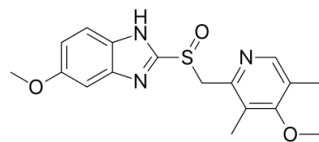
3. Highlight Reaction Centers

<https://metpred.service.pharmb.io/draw/>

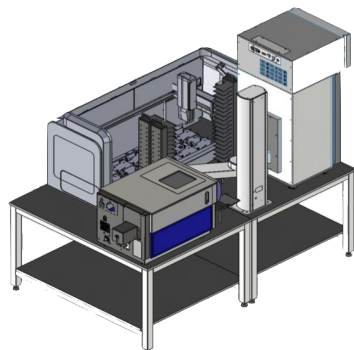
<http://ptp.service.pharmb.io/>



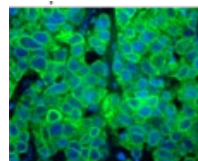
UPPSALA  
UNIVERSITET



chemicals



Robotized lab



images



cell profiles

Chemical DB



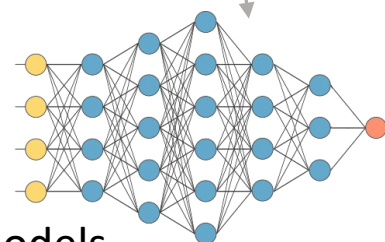
link



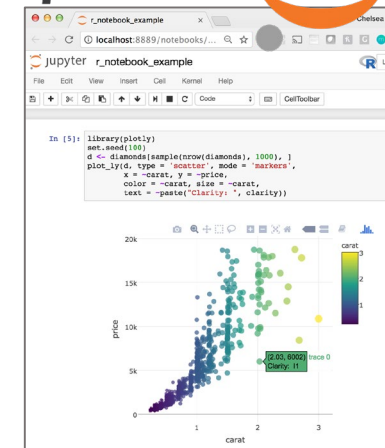
OMERO



CPU/GPU/HPC cloud



Models



Notebooks



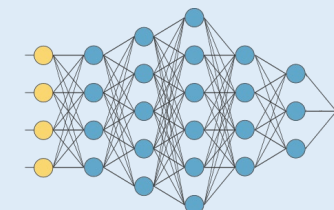
External  
user



Services



Data



Models

Public services



# DevOps and components

- **IaaS**: Infrastructure is portable, testable and simplifies maintenance
  - **Containers**: Components become portable and scalable
  - **Kubernetes**: Resilient container orchestration over multiple nodes
    - logging and monitoring - profiling
    - Hybrid cloud - elasticity
  - **Pachyderm**: Workflows of containers in Kubernetes with data versioning
  - **CICD**: Streamline development and testing
- **DevOps**: Software Developers, Data Engineers and Data Scientists working together in same infrastructure

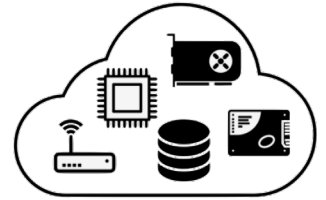


**scaleout**  
SYSTEMS



# Implications: Continuous Analytics

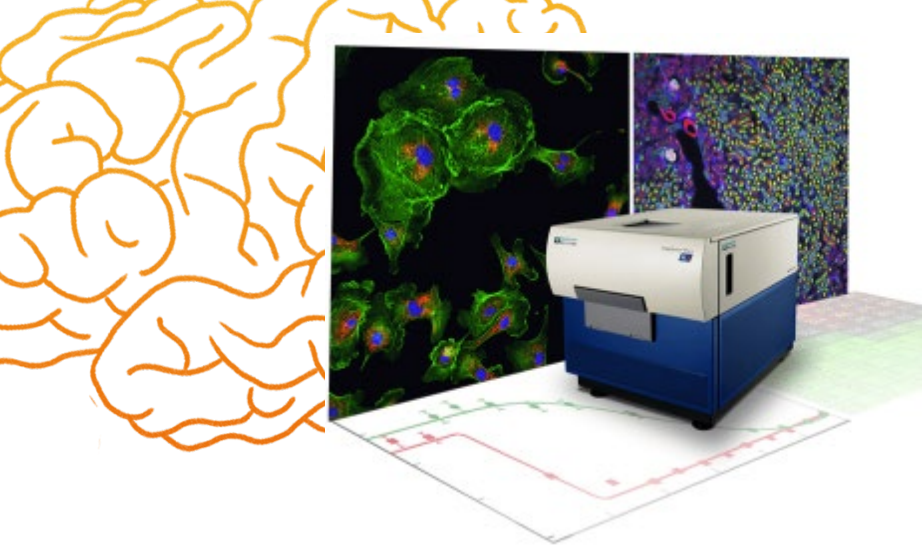
- We can handle the continuous data processing from instruments with robust, resilient data pipelines
- We can continuously re-train models as data is updated
- We can continuously publish data and models
- Agile research group of different competencies
  - Scientists get access to necessary infrastructure
  - DevOps roles, no dedicated sysadmin / tool devel / scientist roles



kubernetes



TensorFlow



Funding:



- Thank you -



Research group website: <http://pharmb.io>

