# ESCAPE to Victory
## Building the infrastructure for the next generation astronomy
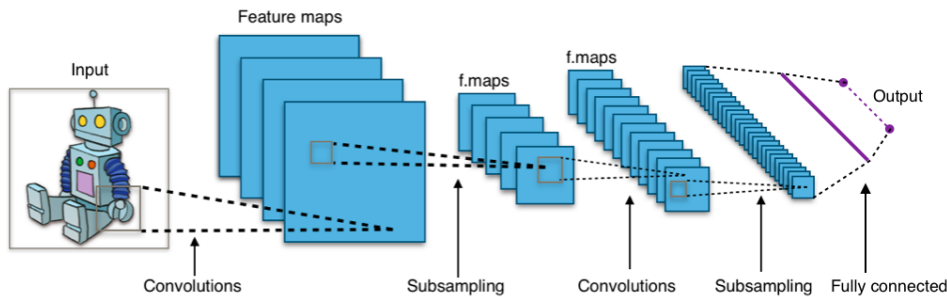
# Challenges
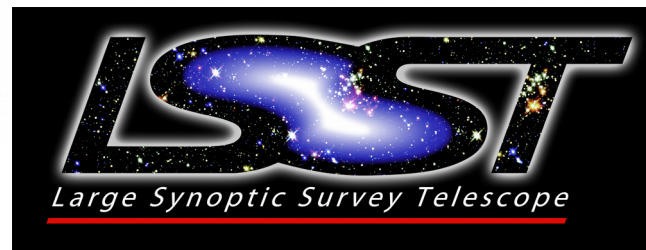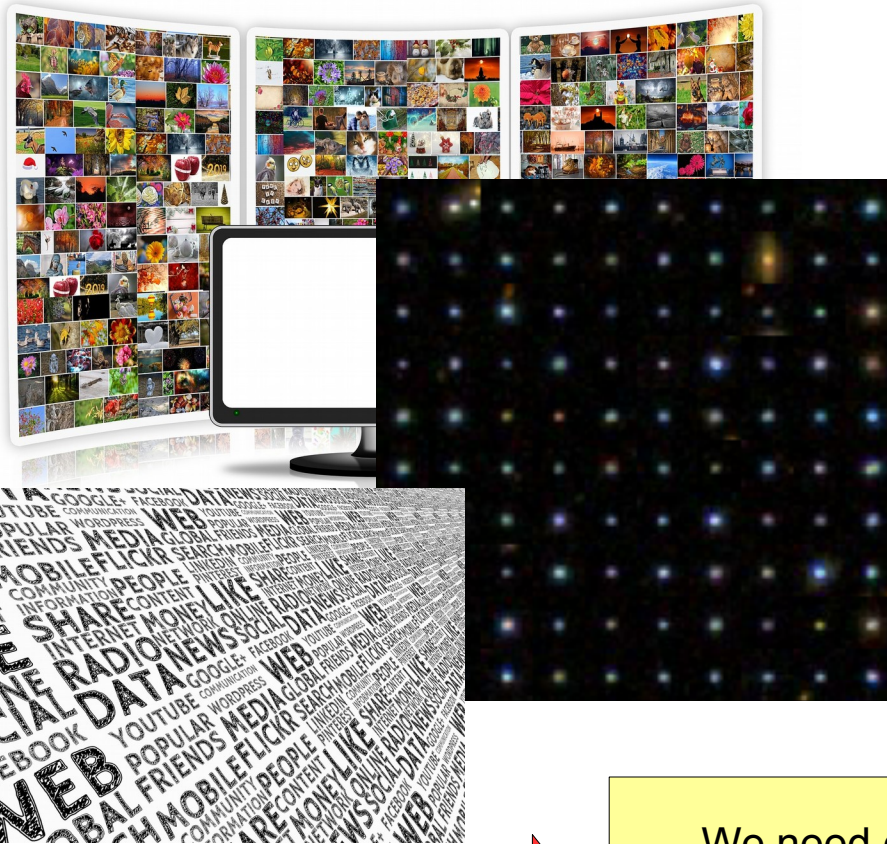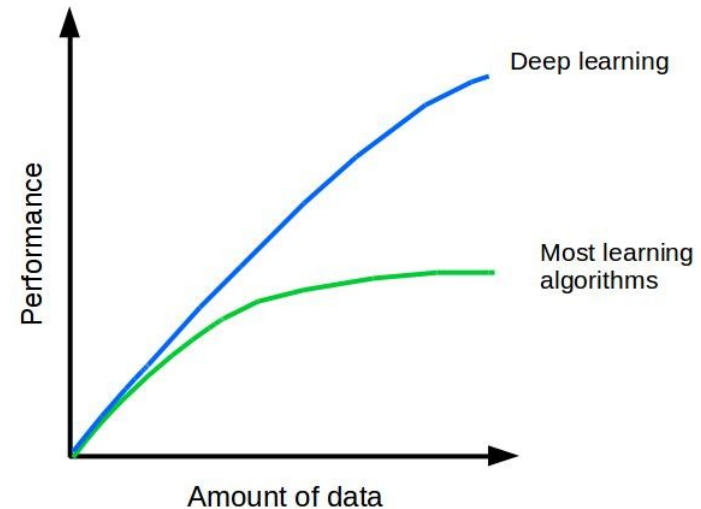
Image from Wikipedia



Big Data by Nick Youngson CC BY-SA 3.0 Alpha Stock Images
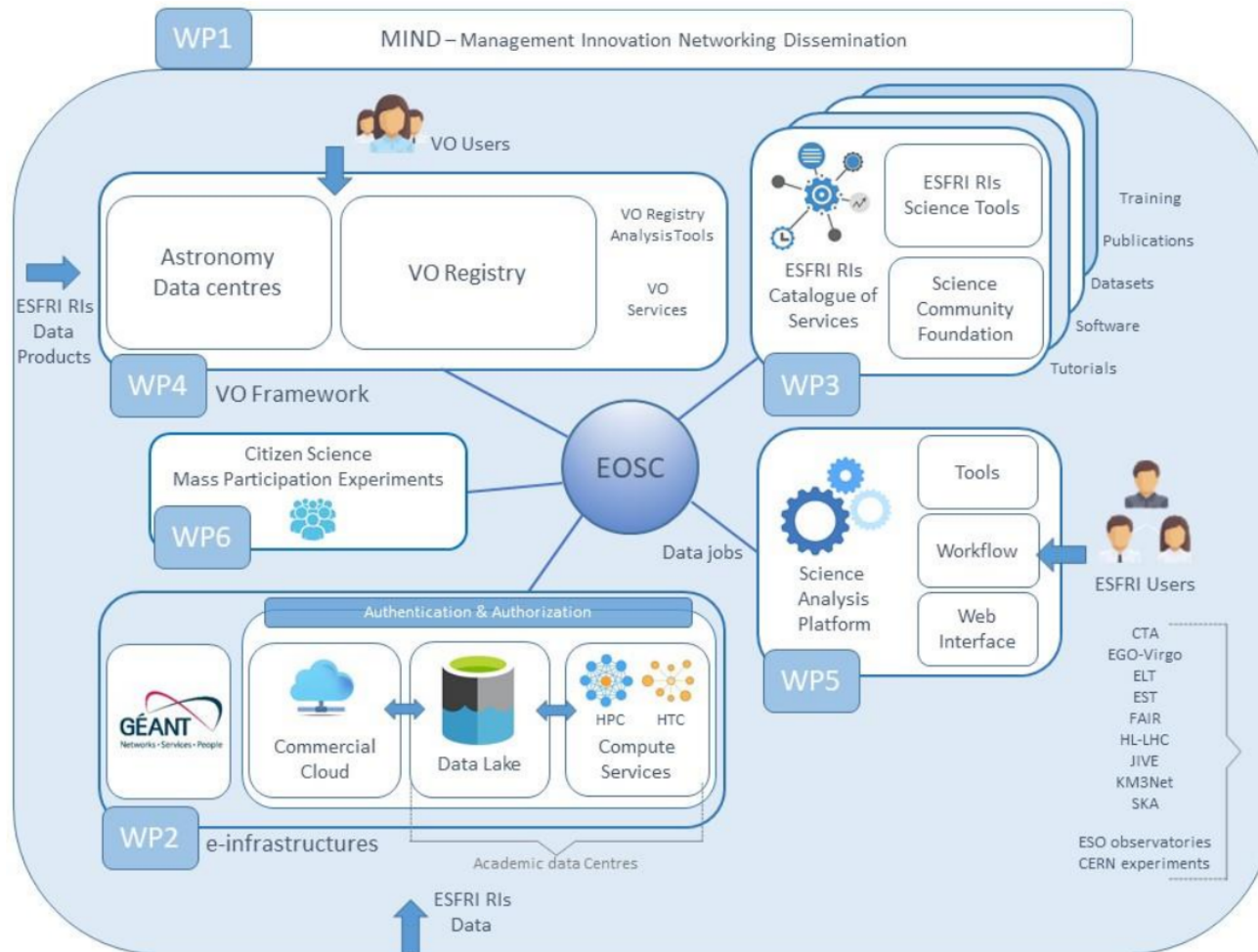
# We are hungry for data!

**BIG DATA & DEEP LEARNING**



We need data,
we need standards,
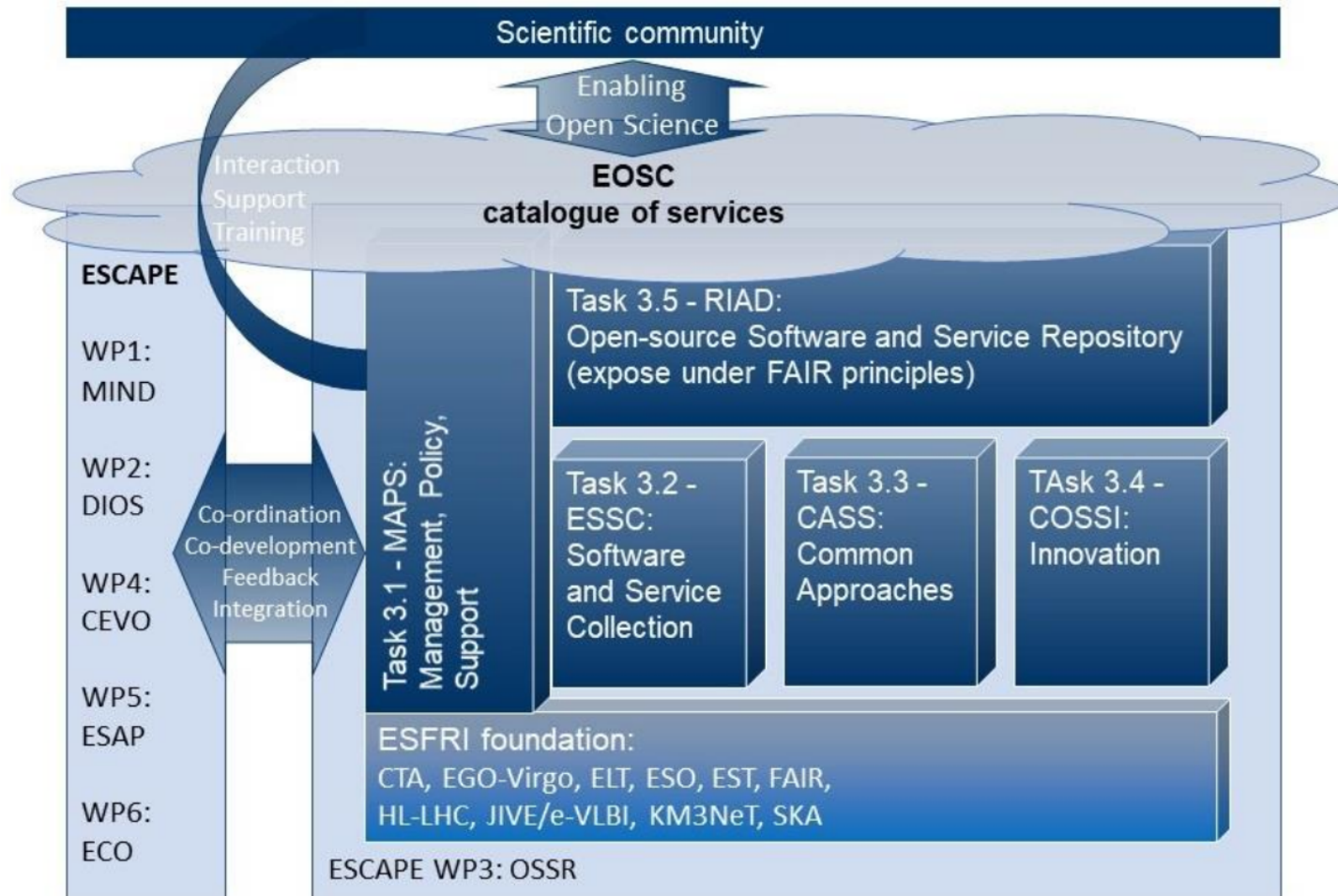we need tools!

# ESCAPE project

## ESCAPE: European Science Cluster of Astronomy & Particle physics ESFRI research infrastructures

- Accessibility to huge amount of data provided by research infrastructures and facilities

- Bring together partners from astronomy and particle physics

- Deliver solutions to ensure integration of data, tools, services and software
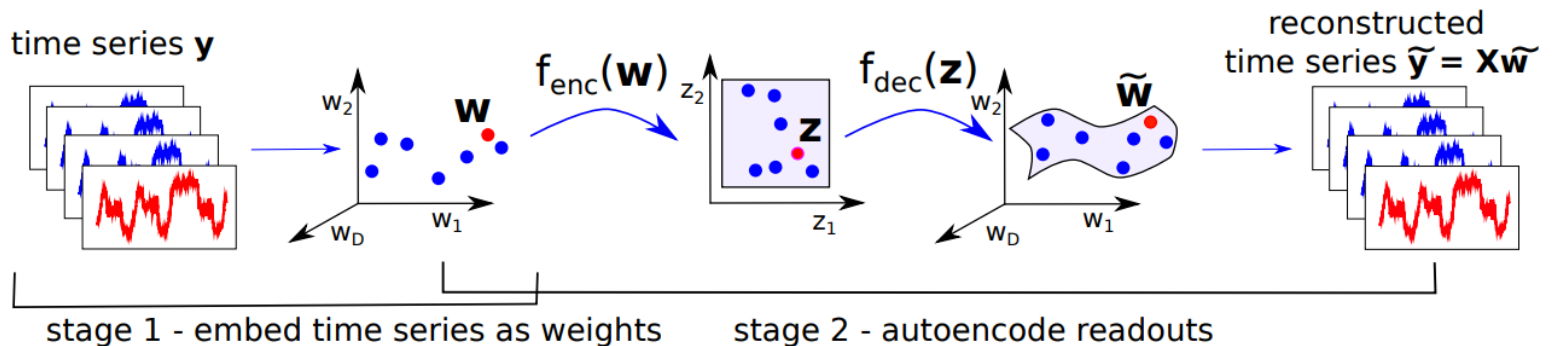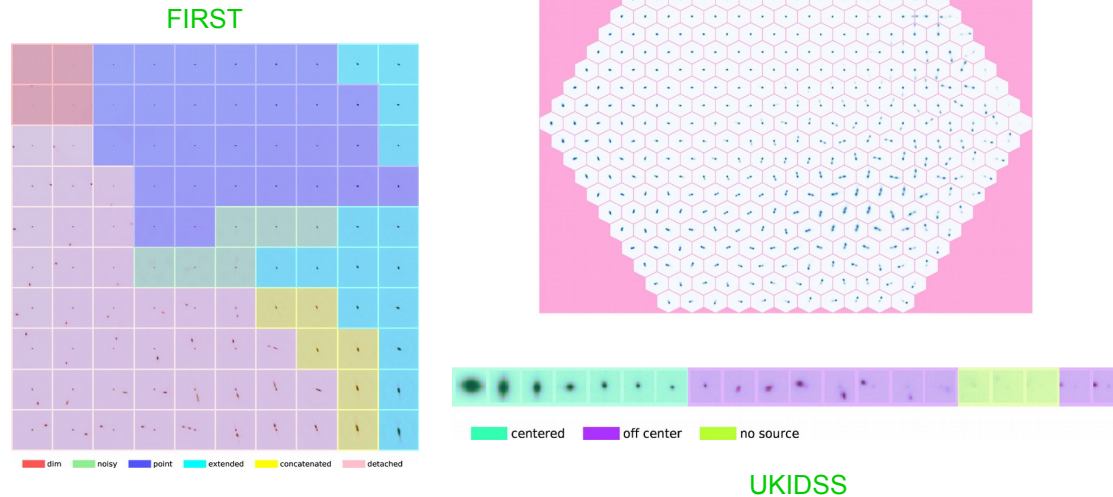
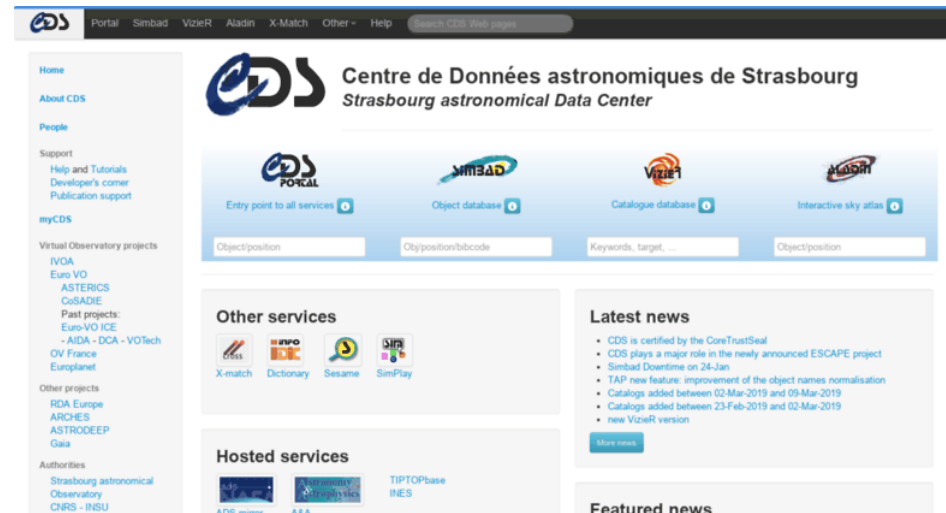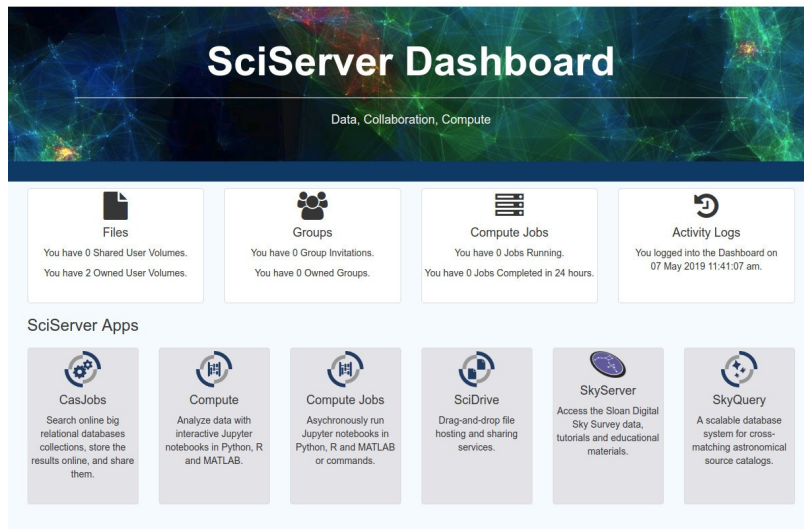- Build standards and ensure interoperability

# Project overview

# WP3

Scientific community

Enabling Open Science

EOSC
catalogue of services

Interaction
Support
Training

**ESCAPE**

WP1:
MIND

WP2:
DIOS

WP4:
CEVO

WP5:
ESAP

WP6:
ECO

Co-ordination
Co-development
Feedback
Integration

Task 3.1 - MAPS:
Management, Policy,
Support

Task 3.5 - RIAD:
Open-source Software and Service Repository
(expose under FAIR principles)

Task 3.2 -
ESSC:
Software
and Service
Collection

Task 3.3 -
CASS:
Common
Approaches

TAsk 3.4 -
COSSI:
Innovation

ESFRI foundation:
CTA, EGO-Virgo, ELT, ESO, EST, FAIR,
HL-LHC, JIVE/e-VLBI, KM3NeT, SKA

ESCAPE WP3: OSSR

# Dimensionality reduction and visualization

## The story so far...

FIRST



dim    noisy    point    extended    concatenated    detached



centered    off center    no source

UKIDSS



time series $y$ → $w_2$, $w_1$, $w_D$, $W$ → $f_{enc}(w)$ → $z_2$, $z_1$, $Z$ → $f_{dec}(z)$ → $w_2$, $w_1$, $w_D$, $\widetilde{w}$ → reconstructed time series $\widetilde{y} = X\widetilde{w}$

stage 1 - embed time series as weights          stage 2 - autoencode readouts

# Challenges: data products

**Heidelberg Institute for Theoretical Studies**
HITS

Working with catalogs is a simple task:



⇨ Problems start with images and spectra!

# Some "simple" tasks...

**Heidelberg Institute for
Theoretical Studies**

**HITS**

1. Given the coordinates, download 28x28 pixel$^2$ images for all the quasars in SDSS.

2. Download some hundreds of thousands of images from FIRST/UKIDSS.

3. Download all the HARPS spectra from ESO archive.

Obtaining data products can be a not easy task

# Task 3

Download all the HARPS spectra from ESO archive:

## Two options

**TAP service**                    **Request system**

# Problems found:

Request system



Return max 10000 rows. All Fields

No file chosen

Limit in the amount of data:
not enough for deep learning!

# Problems found:

TAP service

- Absence of a clear schema browser



- Poor and unclear documentation

# Solutions?

- Python script available on request (thanks Alberto Micol and Martino Romaniello!), but frequent crashes experienced → still investigating the problem, solved by hacking

- Download much slower with respect to request system

**Asking friends and colleagues for data and support
can be helpful, but it is not what standardization and
the ESCAPE paradigm are about!**

# Bringing code to the data

Uploading my code and work on server side could solve many issues…



Who is going to provide and pay for resources?

# First developments

Development of a prototype for:

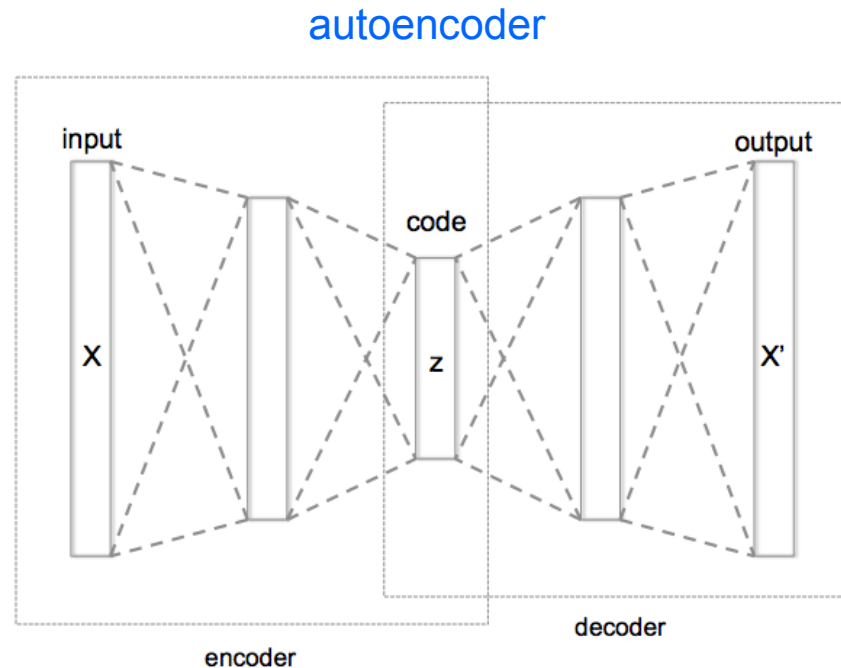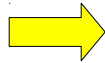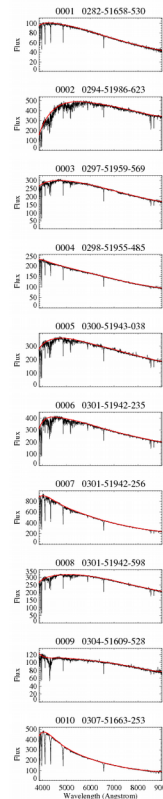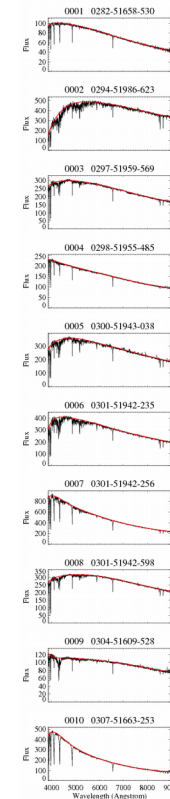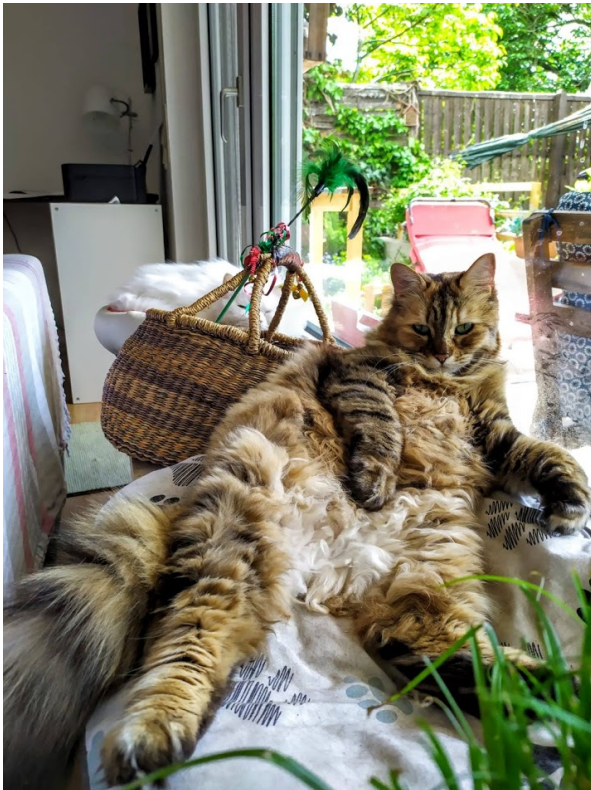Dimensionality reduction and analysis of spectra



autoencoder

representation

# First results

# SEE LIVE DEMO

# Conclusions

**Heidelberg Institute for**
**Theoretical Studies**

HITS

- ESCAPE project is going to be a step to build a new infrastructure for data-intense astronomy

- A lot of work to do:
    - ✗ data products access
    - ✗ building standards
    - ✗ bringing code to the data

- Development of a first prototype → big potential and future integration in web services

- Final question: are we ready for machine learning and big data?

THANKS!
QUESTIONS?