

# **Active Deep Learning In LAMOST DR2 MegaSpectra Archive**

## **The Crucial Role of VO in Modern ML**

**Petr Škoda**

Astronomical Institute of the Czech Academy of Sciences Ondřejov  
Faculty of Information Technology of the Czech Technical University in Prague

**Ondřej Podsztavek**

Faculty of Information Technology of the Czech Technical University in Prague

Supported by OP VVV, Research Center for Informatics,  
CZ.02.1.01/0.0/0.0/16\_019/0000765

IVOA Interoperability Meeting KDD IG  
Paris Observatory, Paris, France  
13-th May 2019

# Thanks to

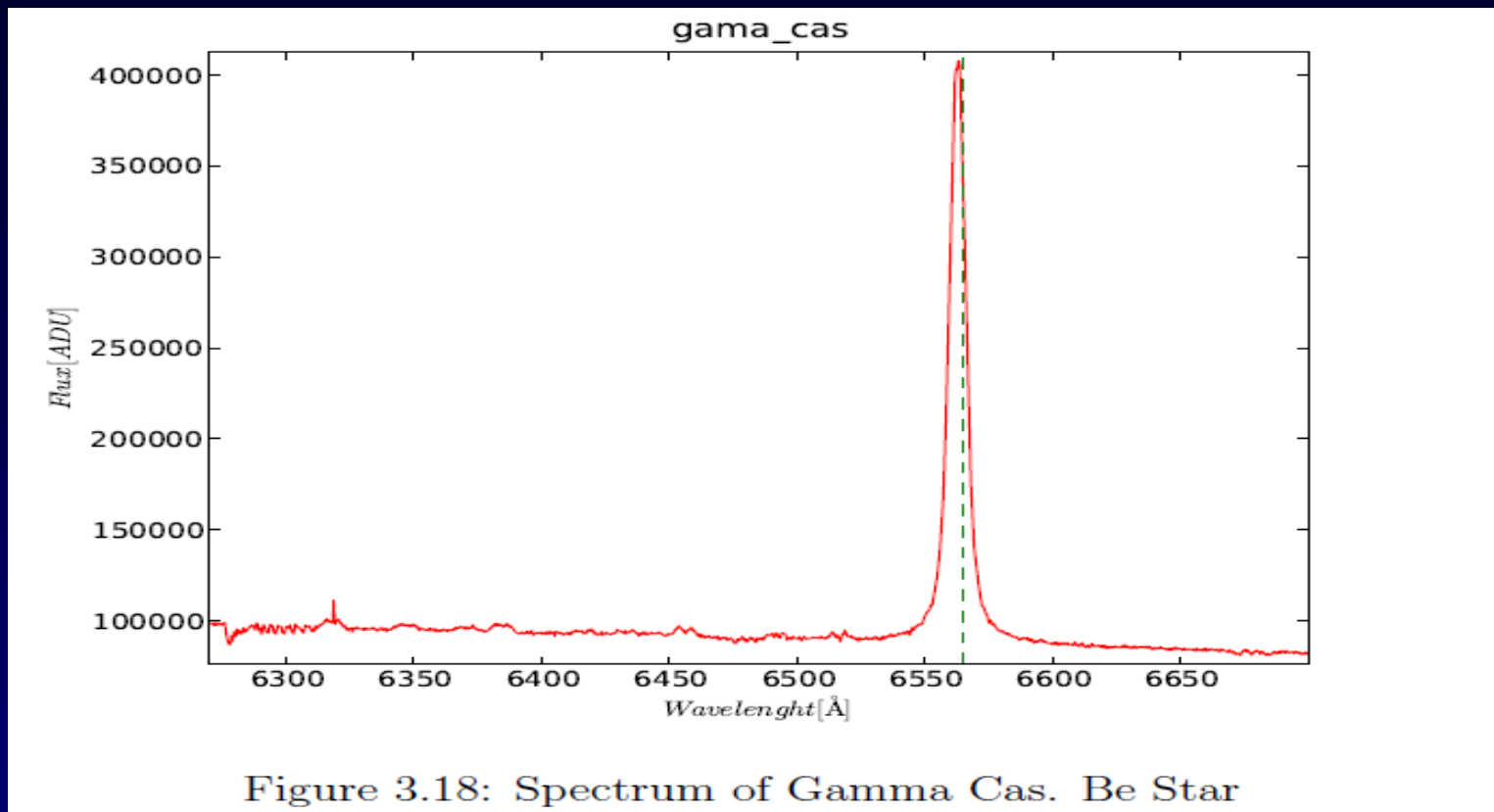
- Astronomical institute CAS- 2m
- LAMOST people (Dongwei Fan, Yue Wu)
- China-VO (Chenzhou Cui)

# Be Stars

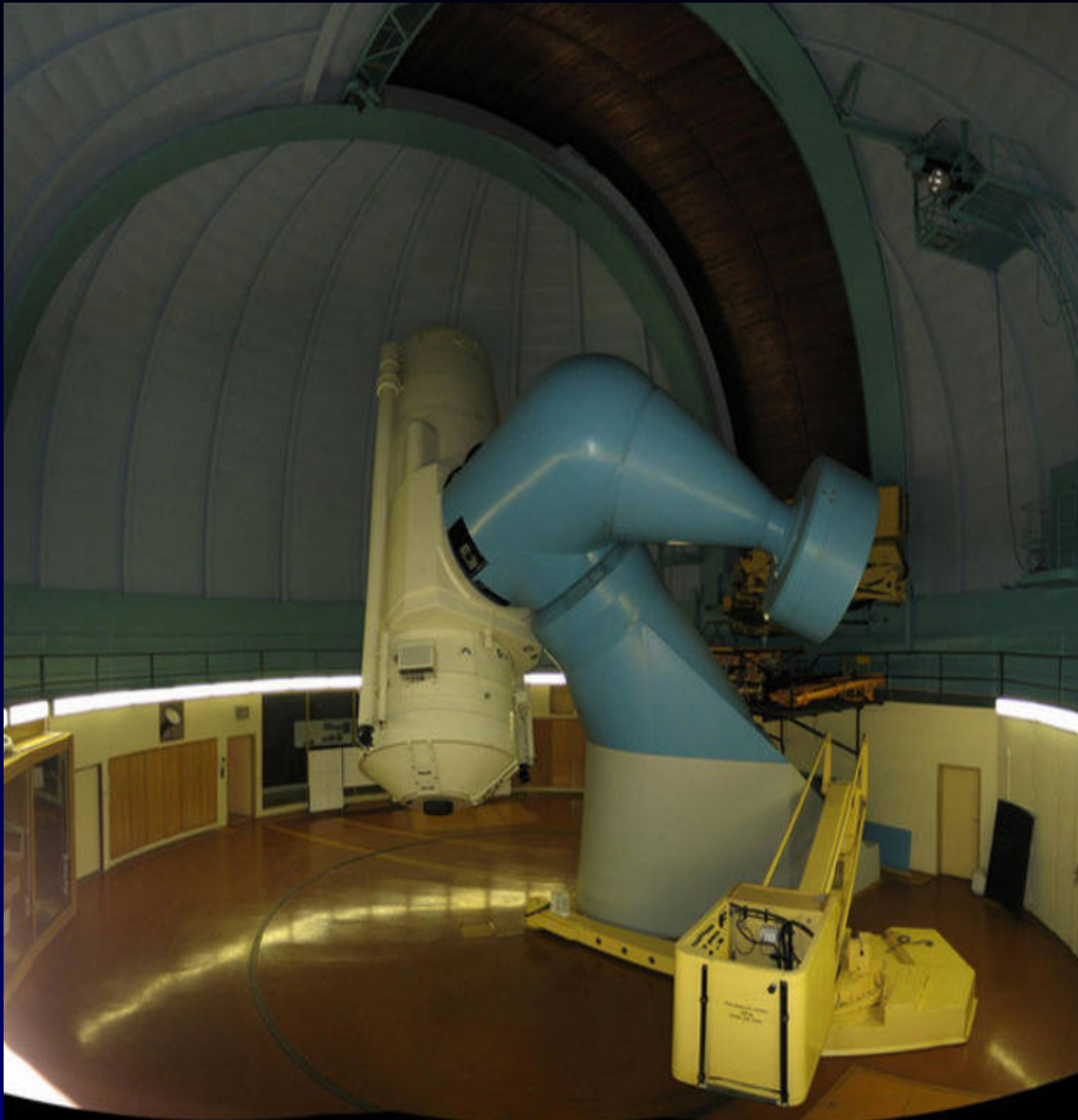
Gamma Cas (Padre Angello Secchi 1866)

Vatican obs predecessor – visual spectrograph

Some have or have had emission in Balmer lines

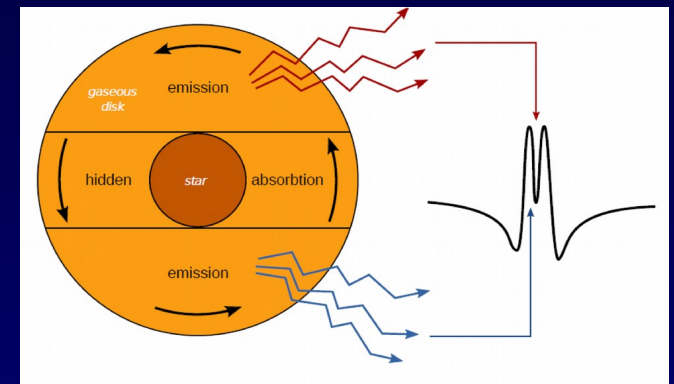


# Ondřejov 2m Perek Telescope (1967)

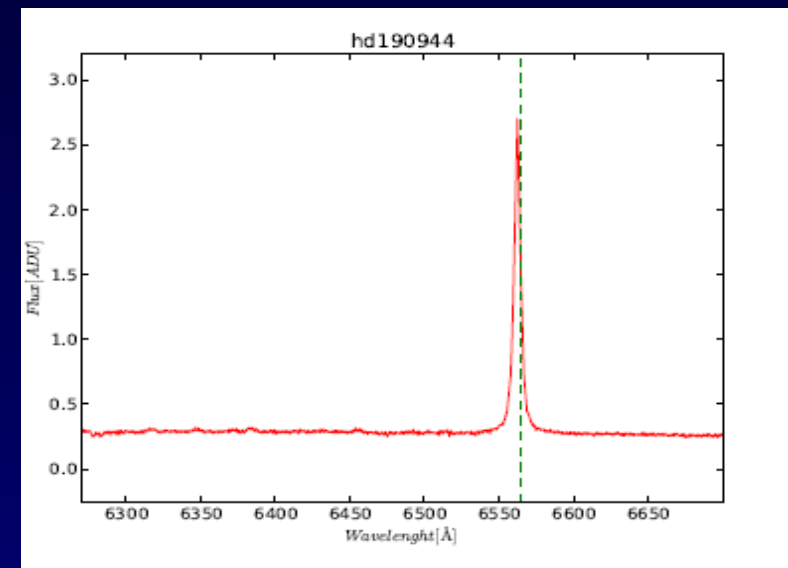
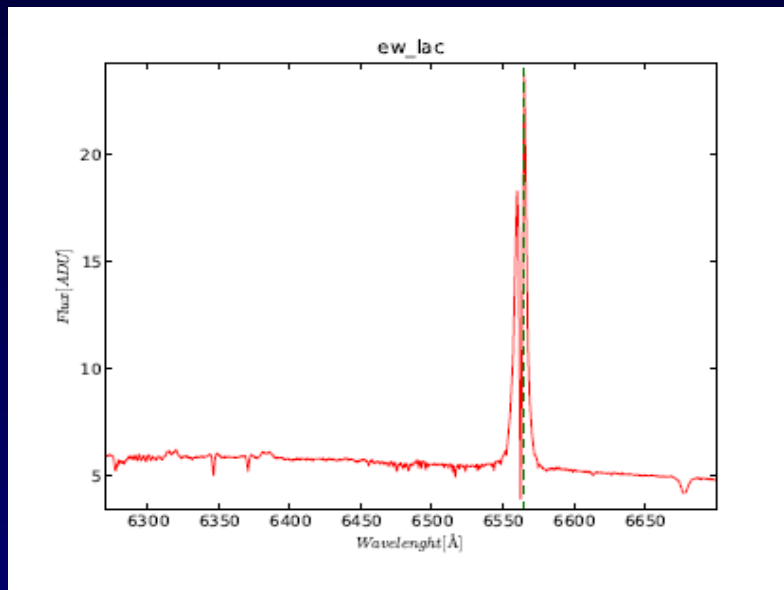
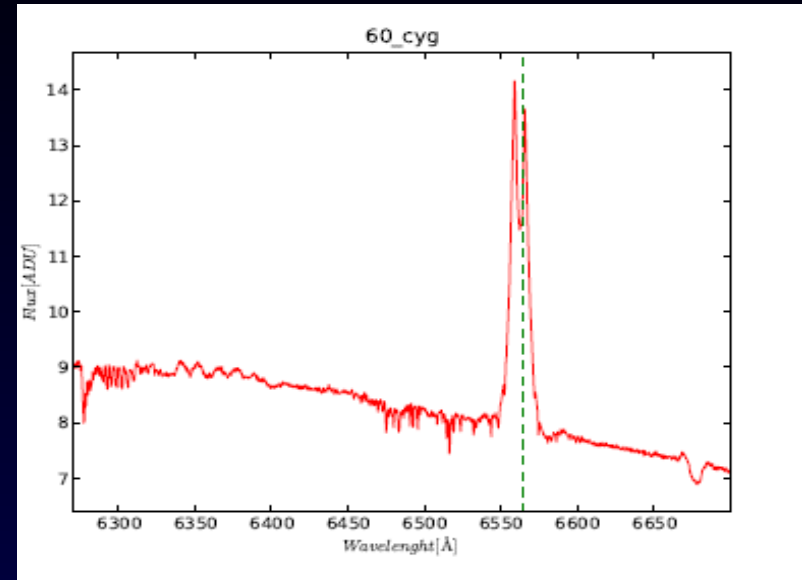
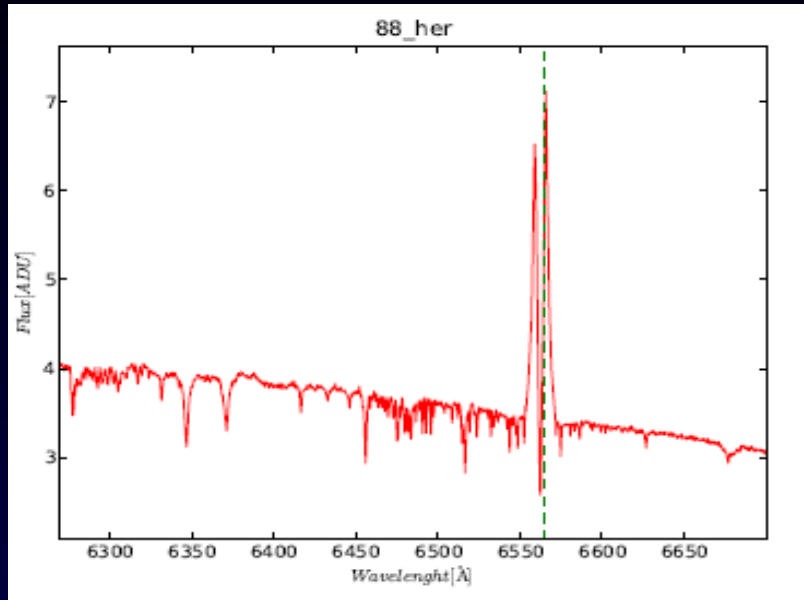


50 years tradition in Be stars  
Archive over 20 000 spectra  
13000 in H $\alpha$  region

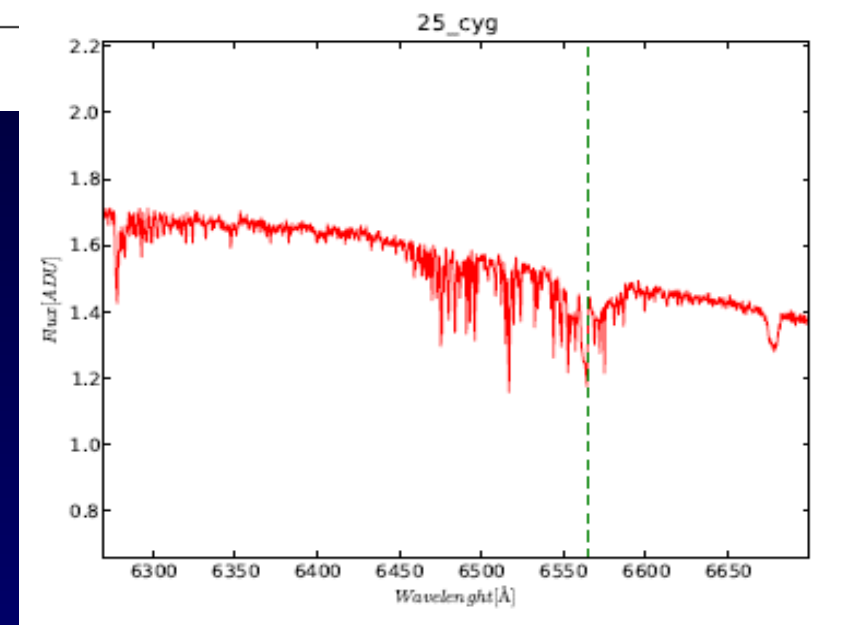
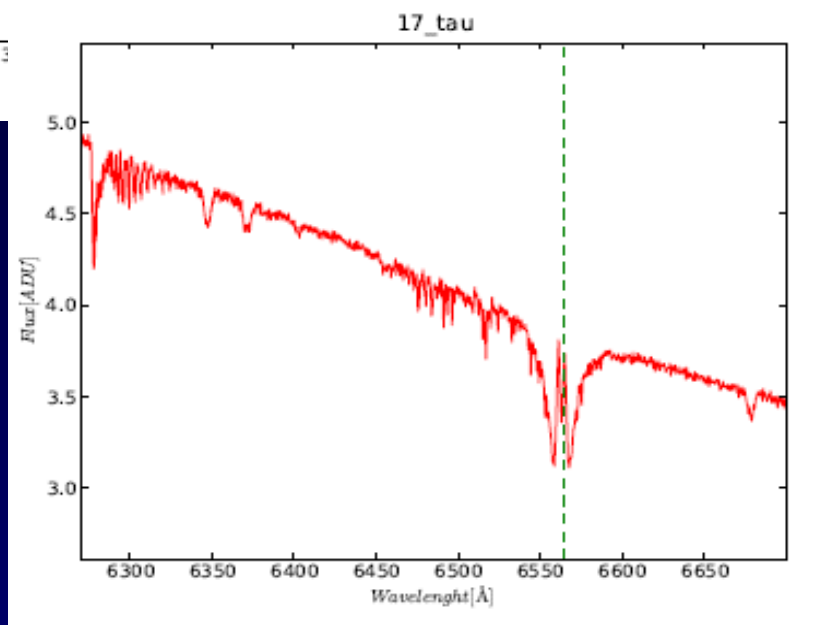
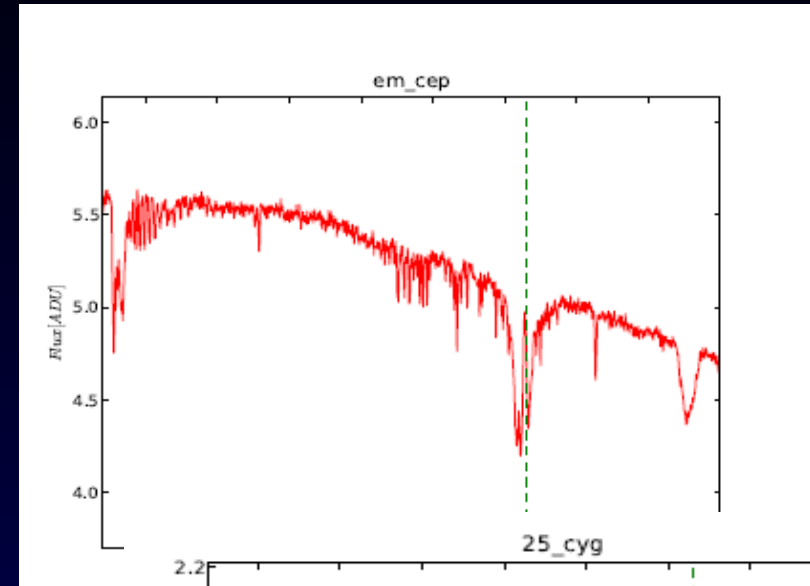
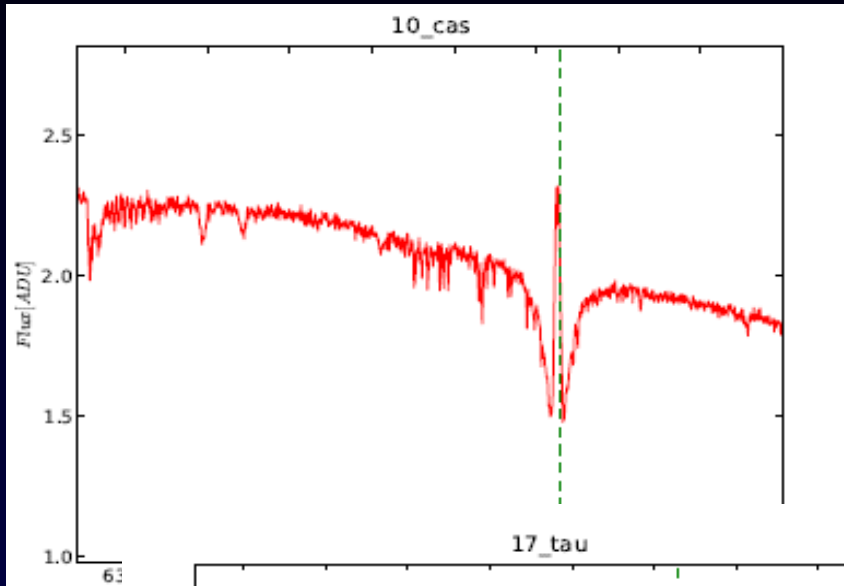
Be mysterious phenomenon  
changes line profile  
episodes of emission  
fast rotate  
Hot (early B types)  
disk or envelope



# Be Stars : Shell lines vs. no absorption

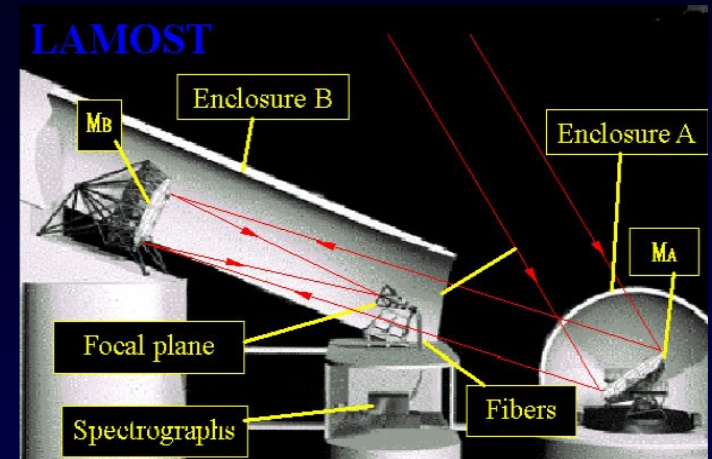
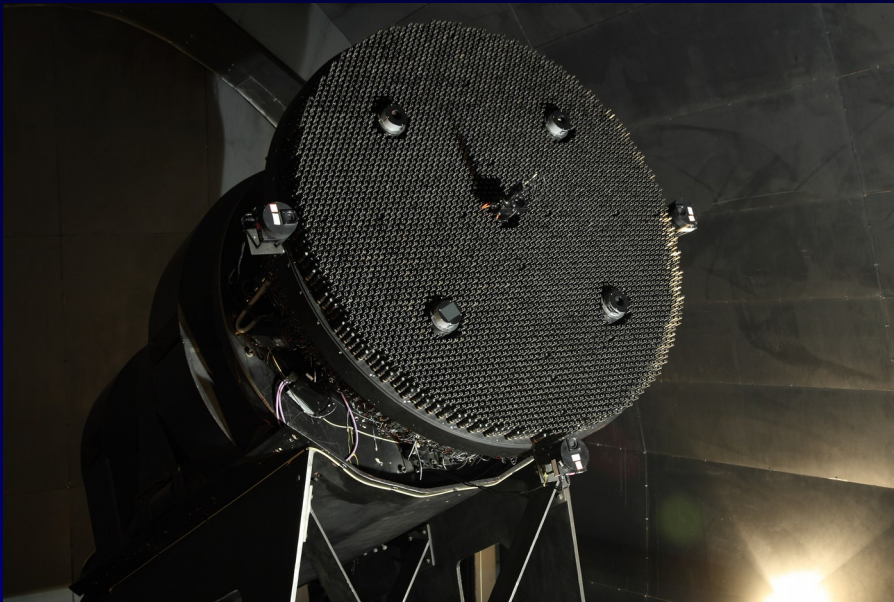


# Be Stars : Emission in absorption



# LAMOST (Guoshoujing)

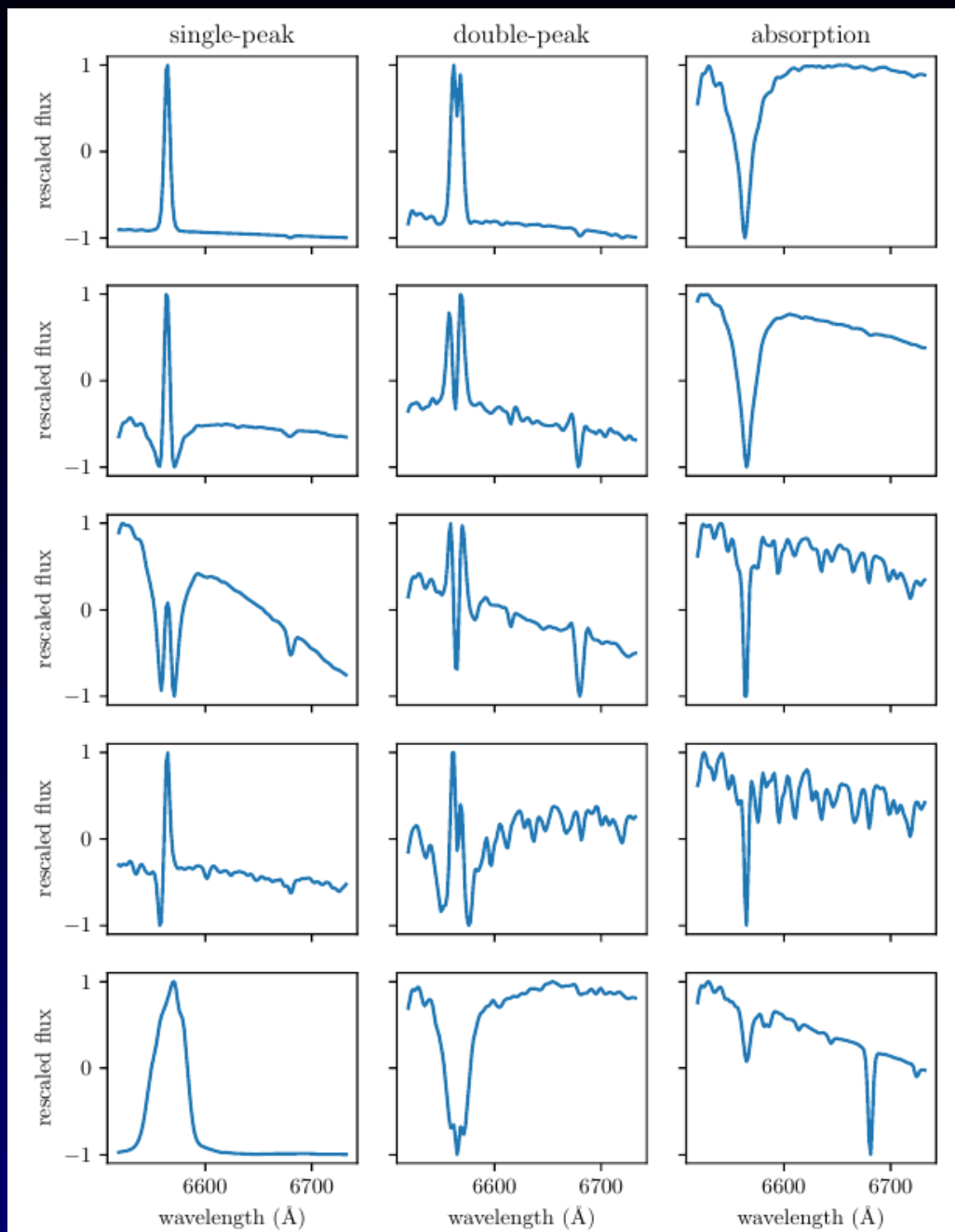
Xinglong- China  
4m mirror (30 deg meridian)  
4000 fibers







# Ondřejov Data Classification



12936 spectra from CCD700

Our TARGET class only Be stars  
=emission or double peak (2+1)

Still not enough labels for DL !

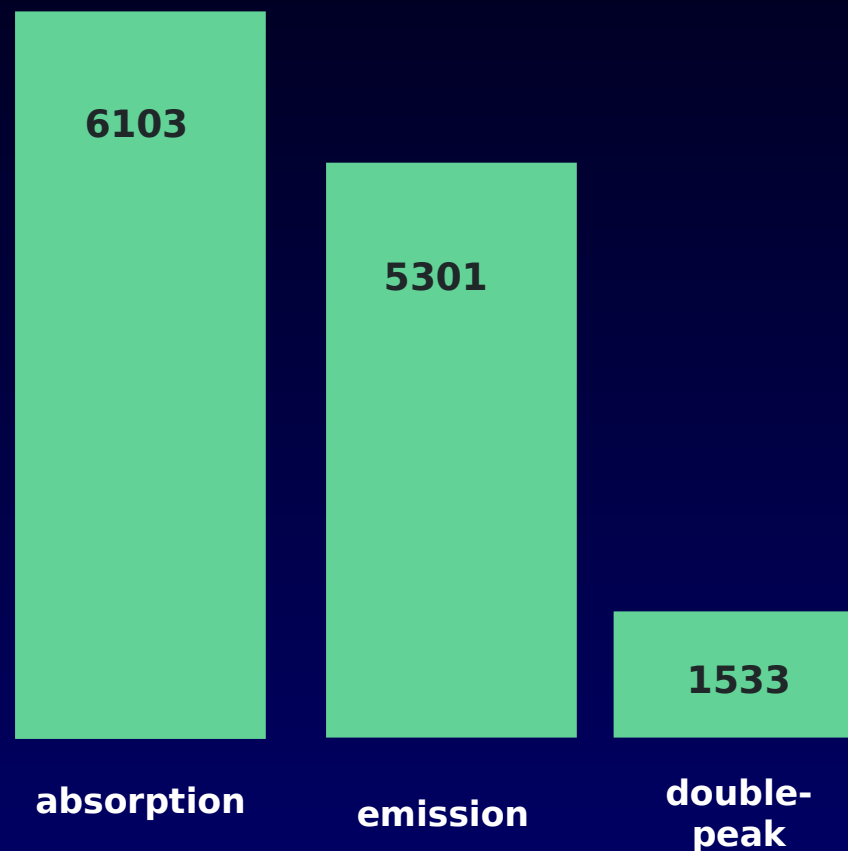
Ondřejov Data Set  
<https://zenodo.org/record/2640971>

# Balancing Classes in Training Set

Synthetic Minority  
Oversampling Technique

SMOTE

Bowyer et al. 2011



# Common Spectra in Both Surveys

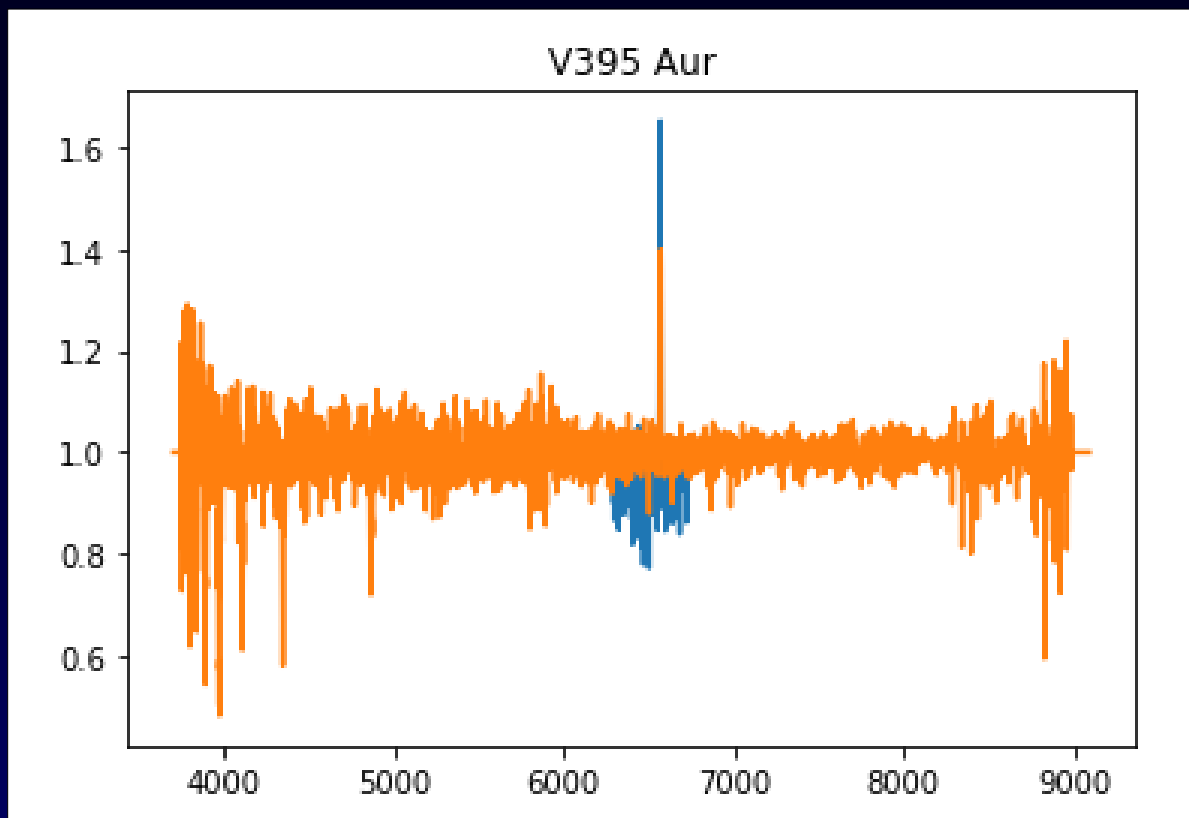
Needed cross-matching in VO  
(Topcat, SPLAT-VO)

LAMOST (4 mil) and  
Ondrejov (13000) on SSAP  
server (DaCHS) + DATALINK

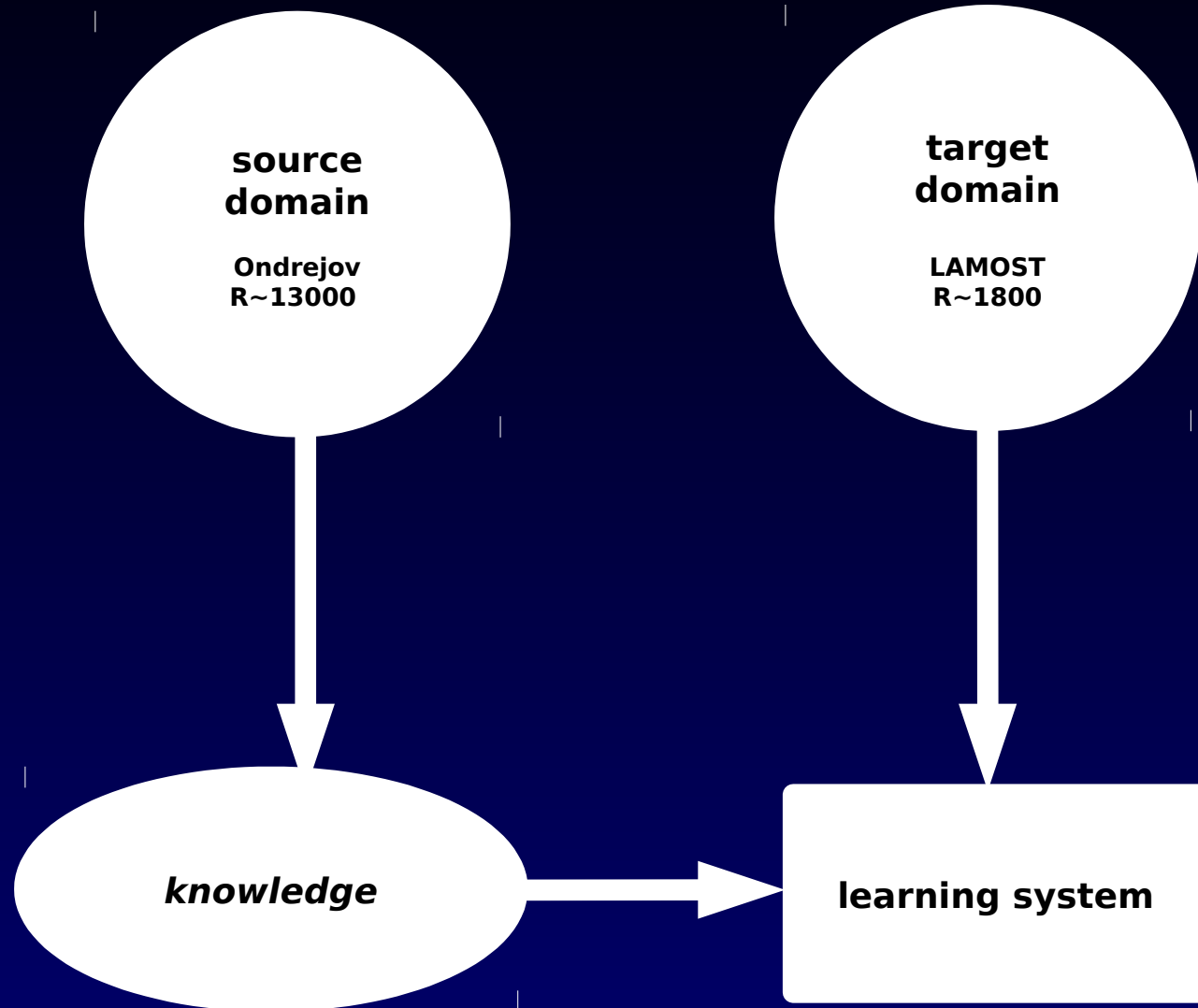
ONLY 22 common (but  
different time of observation =  
changes in profile)

ONLY 4 emission line

No way to proper  
training set!

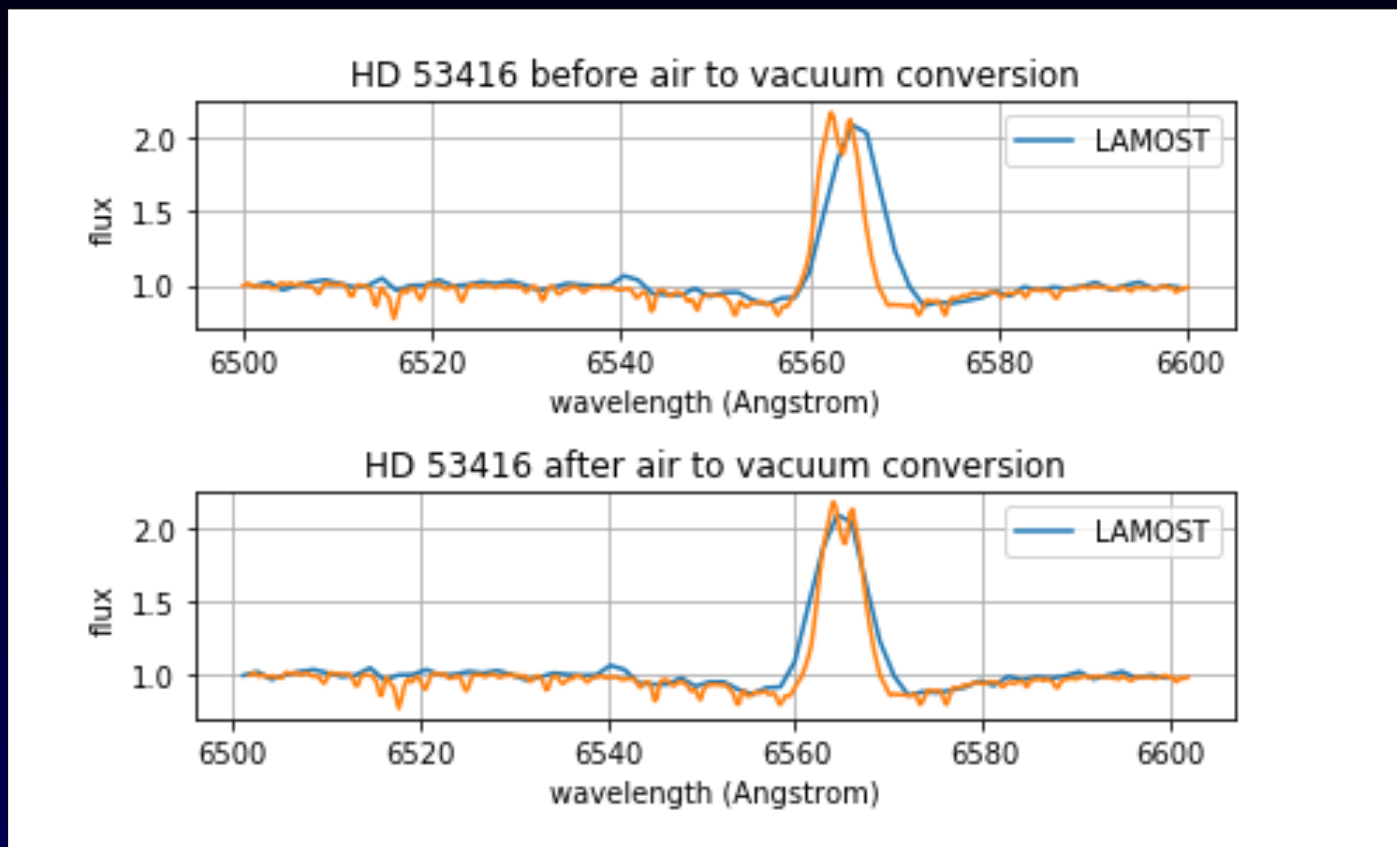


# Domain Adaptation



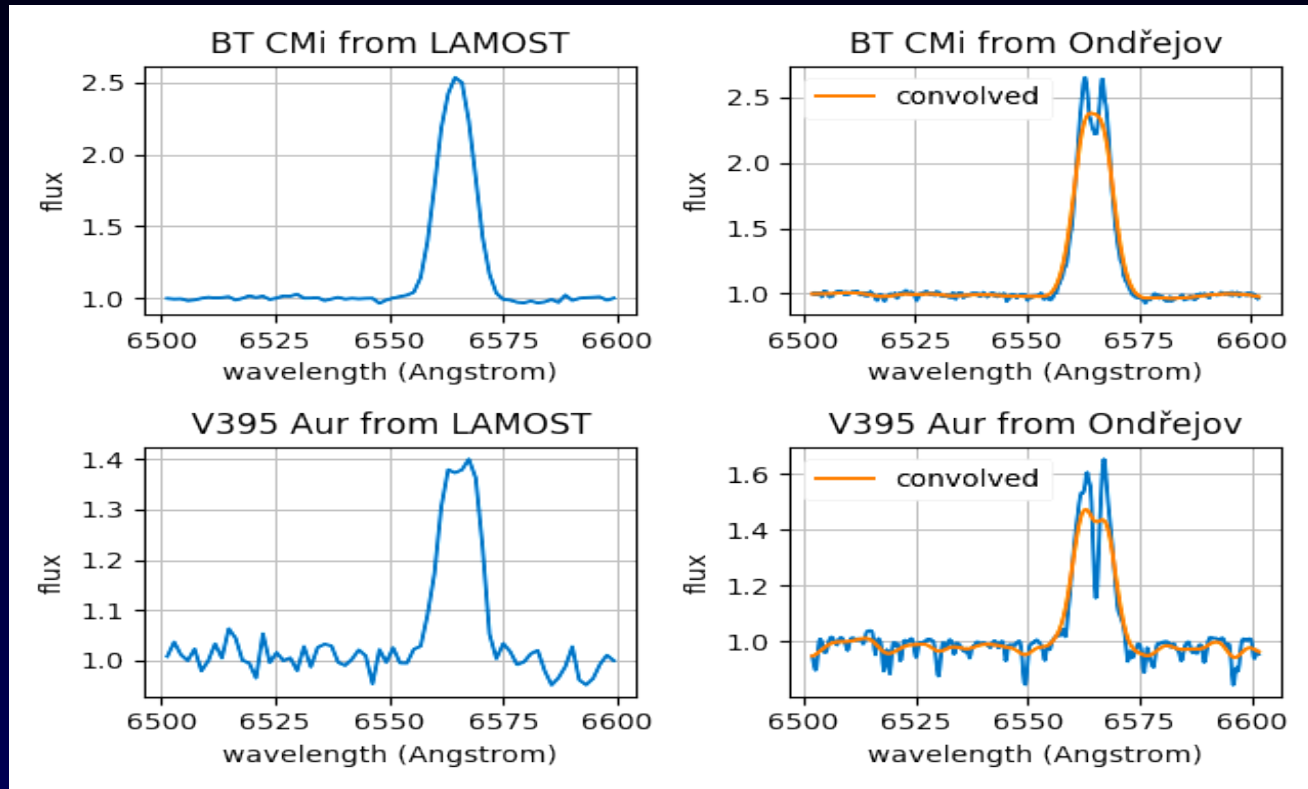
Transfer Learning

# Domain Adaptation



Wavelength transform from air to vacuum

# Domain Adaptation



Change of spectral resolution by **Gaussian blur** (13000  $\rightarrow$  1800)  
**Rebinning** to grid of 140 pixels in 6519 to 6732 A  $\rightarrow$  FEATURE VECTOR

# Normalisation od Data

Continuum normalisation – rectification (continuum on 1.0, limit at 0)

Ondřejov – automatic pipeline

*experiments show that (at least for our Halpha region ) the continuum normalisation is not needed (training on unrectified gives same validation accuracy 0.96 )*

Rescaling to zero mean, unit variance - common in ML

# Deep Convolutional Net

deep network to have representation power  
and **dropout** to reduce overfitting.

Inspired by **VGGNet**  
adapted to 1D spectrum.

No feature extraction!

Training:  
TENSORFLOW+KERAS on GTX980 GPU

input (140 pixel spectrum)
conv3-64
conv3-64
maxpool2
conv3-128
conv3-128
maxpool2
conv3-256
conv3-256
maxpool2
fc-512
fc-512
softmax



# Active Learning

Human (Oraculum) involved

Random/**Uncertainty** Sampling:

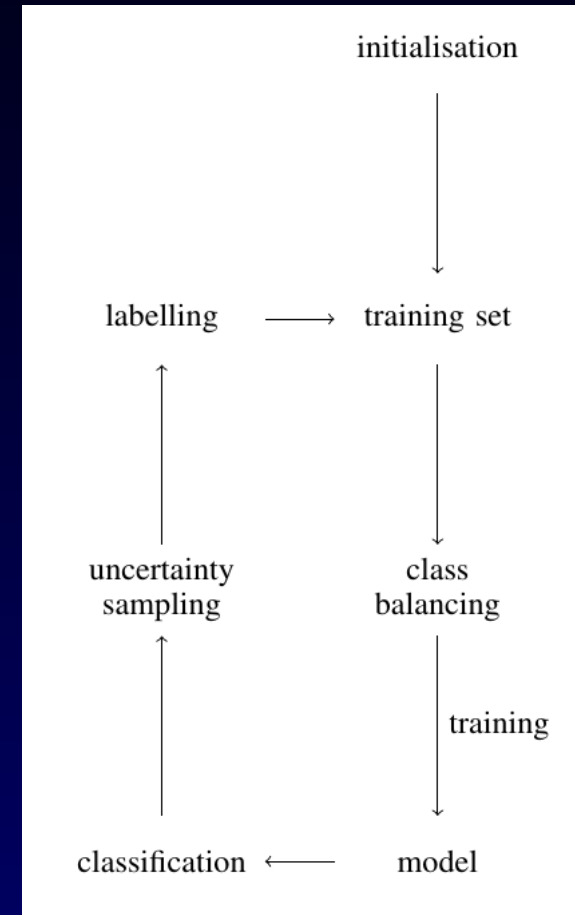
From predicted **TARGET** class (single or double peak) selected 100

randomly/**with highest entropy**

**Visual check** : re-classification  
(confirm, change, put in uninteresting)

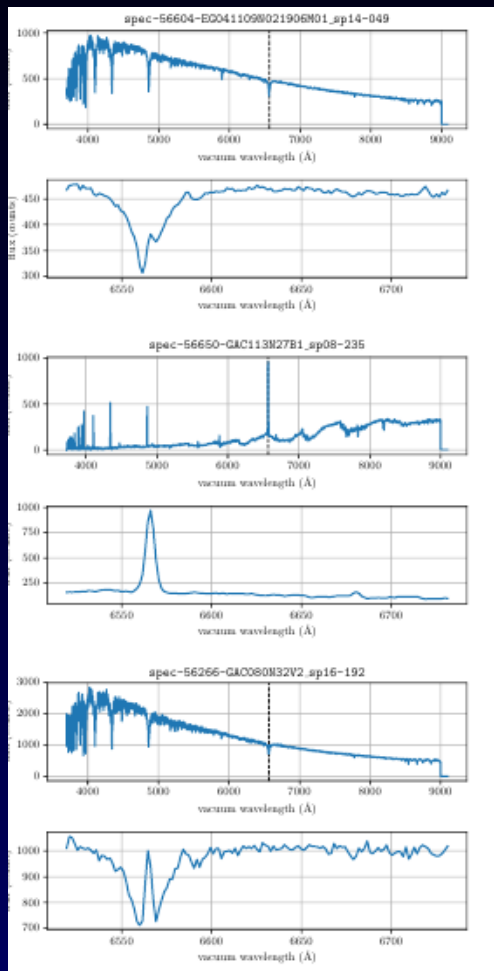
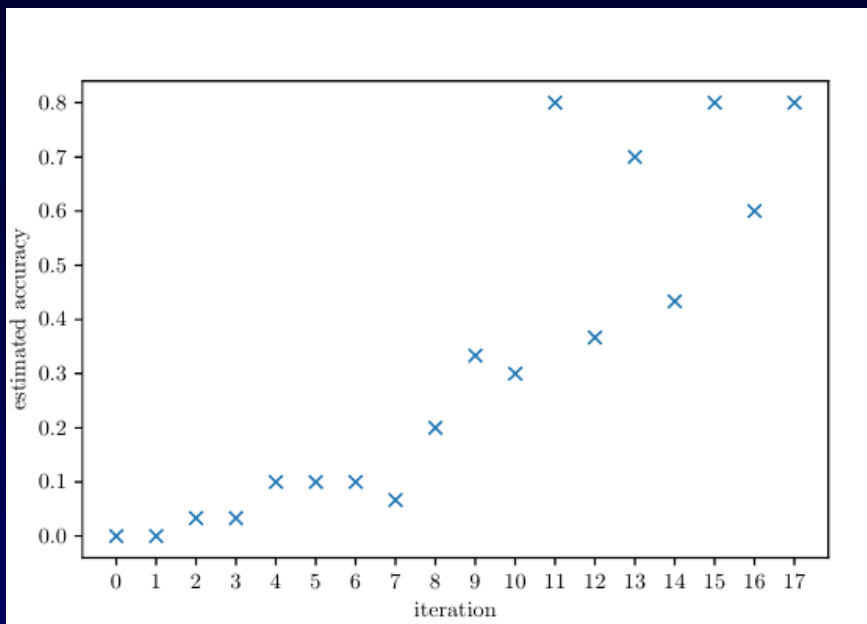
**These data** added to training set

Repeat  
until few misclassifications  
(16 times)

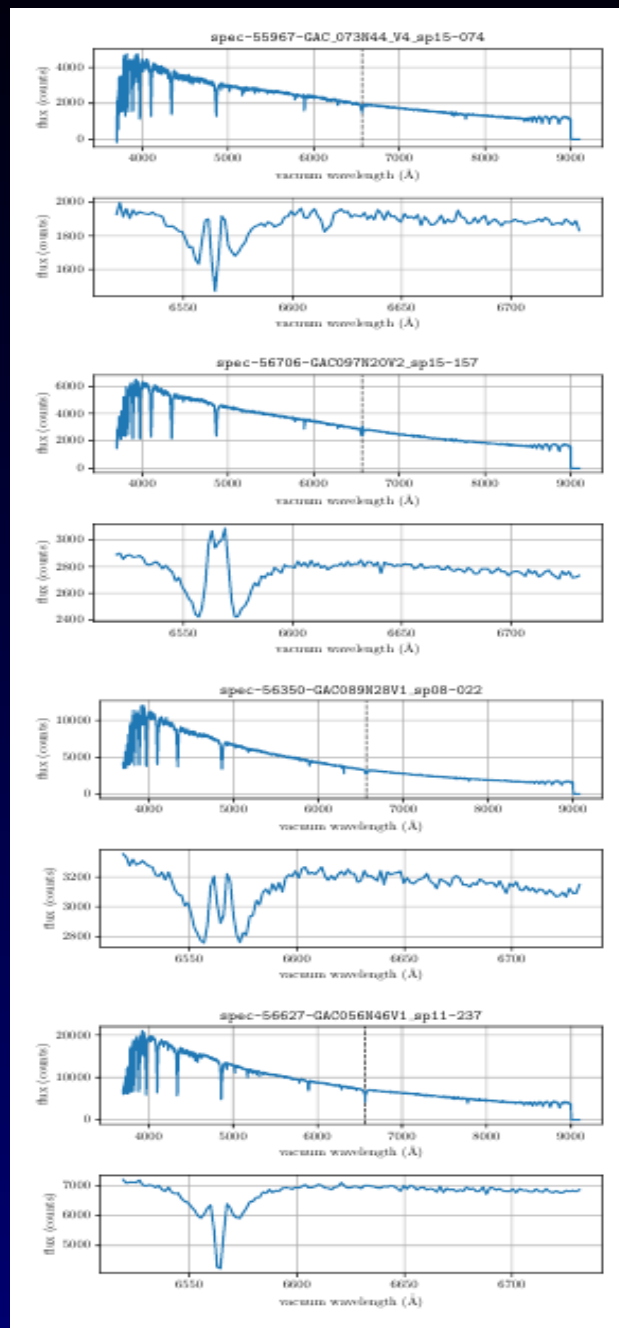


# Active Learning (< 6 % error)

Predicted class	single-peak	Actual class double-peak	uninteresting
single-peak	97.6% (3 641)	1.4% (53)	1.0% (37)
double-peak	2.8% (18)	94.0% (609)	3.2% (21)



Single peak



Double peak

# Results

Number of SPECTRA (multiple)

4379 candidates (from 4 mil)

58 bad – but still interesting – e.g. LAMOST HVS-1

3731 single peak

648 double peak (but complex shapes)

New objects – 664 in SIMBAD (Be, CV, Seyfert Gal...)

1013 (948 objects) NEW , NOT known

Visual check , XMATCH in VO, DSS2 , SDSS in Aladin

Most are correct - ? physical origin (YSO, CV, M, Novae)

Unreliable classification from LAMOST 1D pipeline (F, A/B)

A lot of new Be stars

**Many SUPRISES !!!**

# Comparison with Other Method

RAA 2016 Vol. 16 No. 9, 138 (12pp) doi: 10.1088/1674-4527/16/9/138  
<http://www.raa-journal.org> <http://iopscience.iop.org/raa>

*Research in  
Astronomy and  
Astrophysics*

## A catalog of early-type emission-line stars and H $\alpha$ line profiles from LAMOST DR2

Wen Hou<sup>1,2</sup>, A-Li Luo<sup>1,2</sup>, Jing-Yao Hu<sup>1</sup>, Hai-Feng Yang<sup>1,2,3</sup>, Chang-De Du<sup>1,2</sup>, Chao Liu<sup>1</sup>, Chien-De Lee<sup>4</sup>, Chien-Cheng Lin<sup>5</sup>, Yue-Fei Wang<sup>6</sup>, Yong Zhang<sup>6</sup>, Zi-Huang Cao<sup>1</sup> and Yong-Hui Hou<sup>6</sup>

<sup>1</sup> Key Laboratory of Optical Astronomy, National Astronomical Observatories, Chinese Academy of Sciences, Beijing 100012, China; [lal@bao.ac.cn](mailto:lal@bao.ac.cn), [whou@bao.ac.cn](mailto:whou@bao.ac.cn)

<sup>2</sup> University of Chinese Academy of Sciences, Beijing 100049, China

<sup>3</sup> School of Computer Science and Technology, Taiyuan University of Science and Technology, Taiyuan 030024, China

<sup>4</sup> Institute of Astronomy, National Central University, Jhongli

<sup>5</sup> Shanghai Astronomical Observatory, Chinese Academy of Sciences, Shanghai 200030, China

<sup>6</sup> Nanjing Institute of Astronomical Optics & Technology, National Astronomical Observatories, Chinese Academy of Sciences, Nanjing 210042, China

$$\sum_{i=-5}^5 f_{\text{obs}}[n_0 + i] / 11 > f_{\text{conti}}[n_0],$$

$$\sum_{i=-1}^1 f[n_0 + i] / 3 > \sum_{i=-2}^2 f[n_0 + i] / 5$$

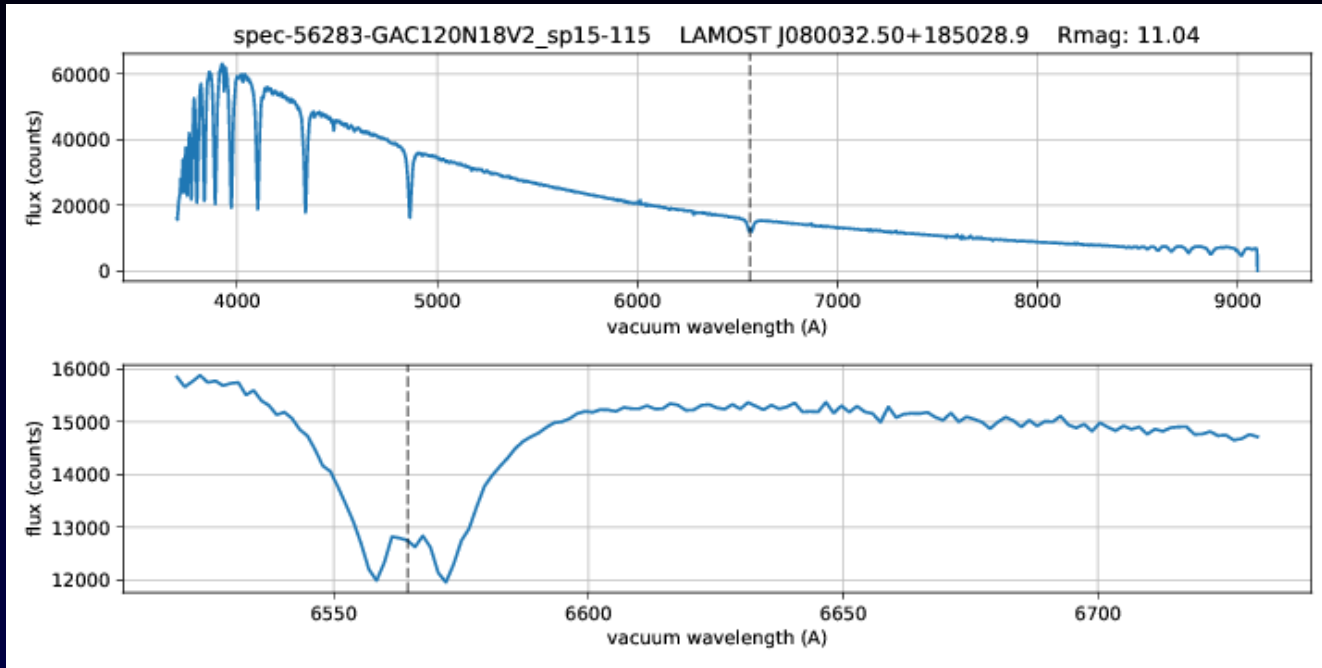
and

$$\max \left( f_{\text{obs}}[n_0 - 1 : n_0 + 1] \right) \geq \max \left( f_{\text{obs}}[n_0 - 2 : n_0 + 2] \right)$$

## Integral pixel statistics on different intervals around H $\alpha$

Hou (2016) DR2 - catalog of 11205 emission stars  
we have 2644 of them - but 11k not well justified !  
(VO xmatch, SPLAT visualization ....)

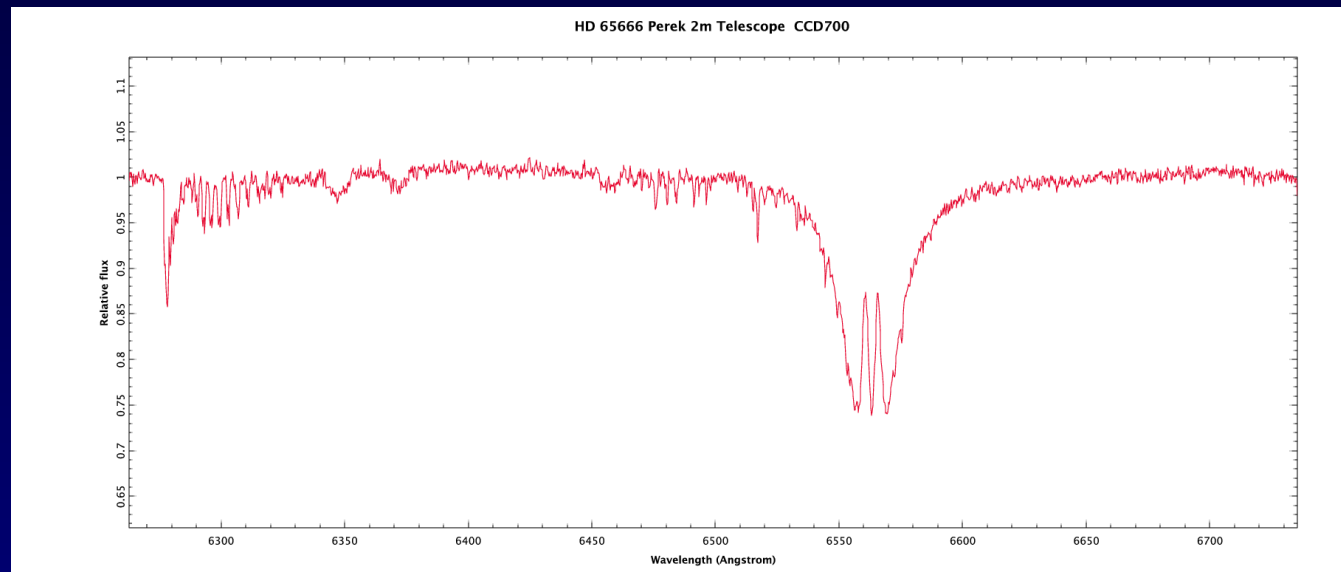
# Results - Unknown Bright Be star



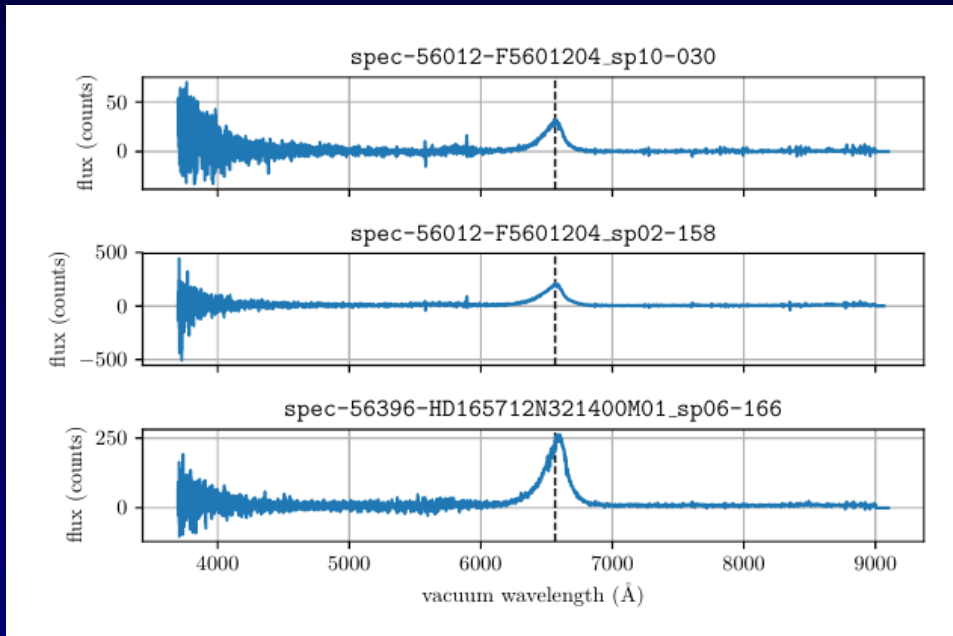
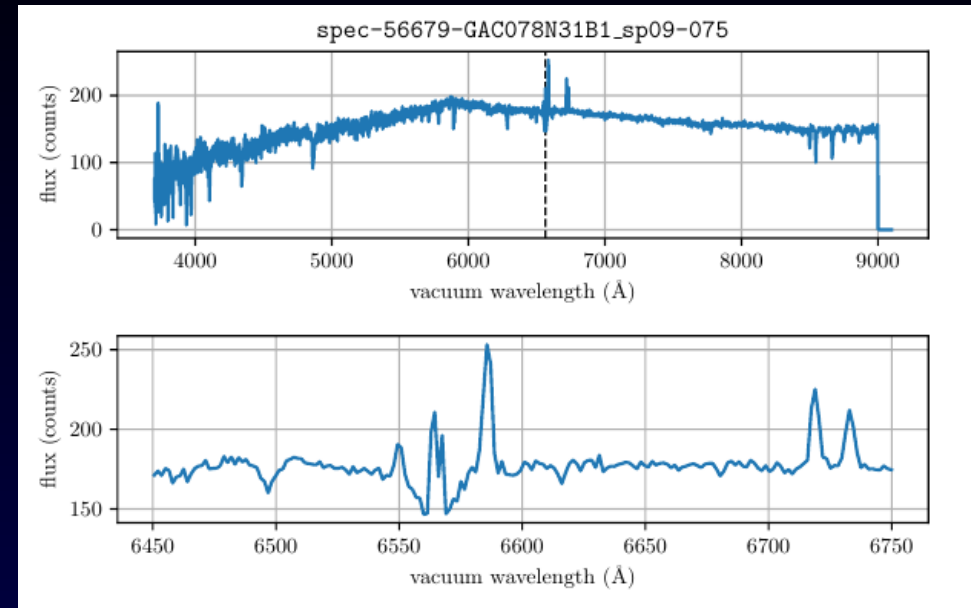
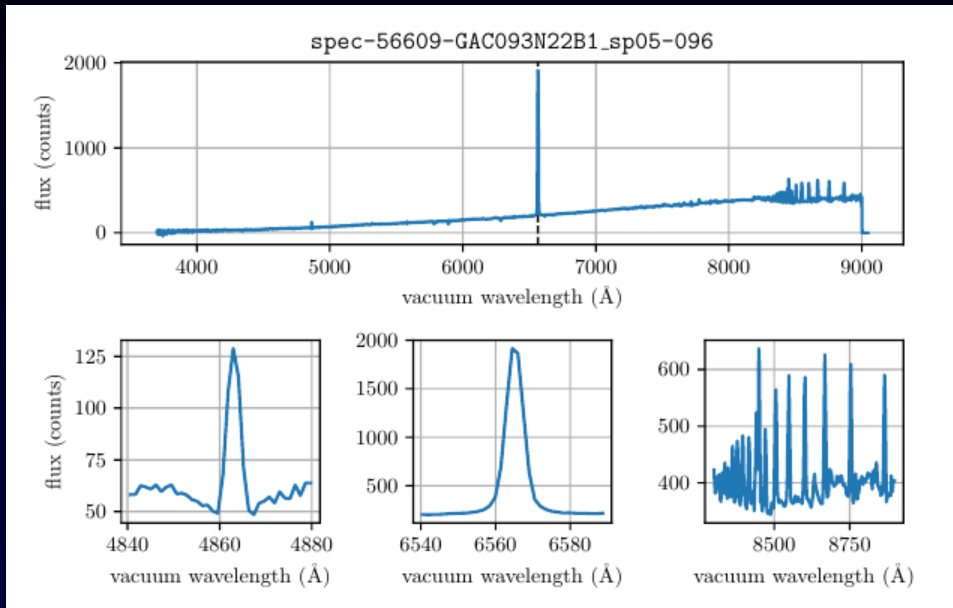
LAMOST

6.5 arcsec distant star 7 mag

CCD700 OND



# Results - Interesting



Supernovae ???

In DSS2 /SDSS there are galaxies in 10arcsec around .....

# Confusion in Unique Identification

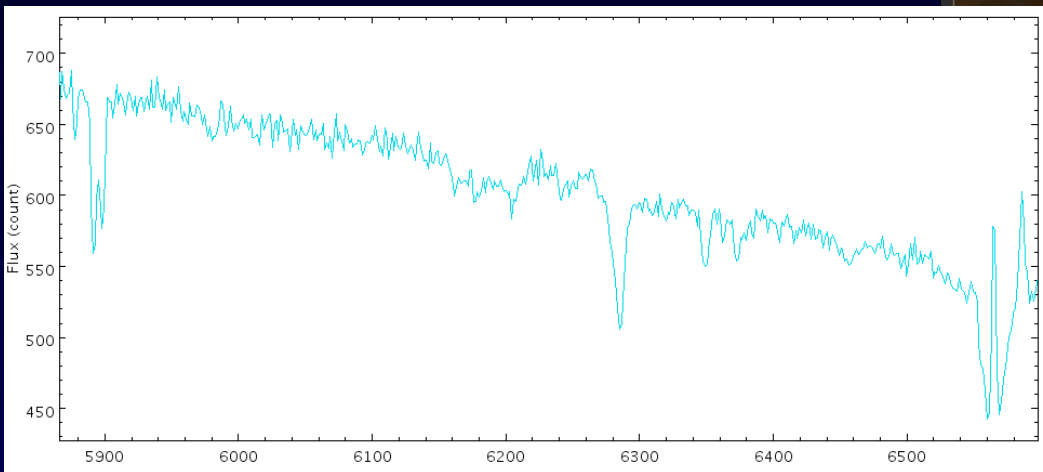
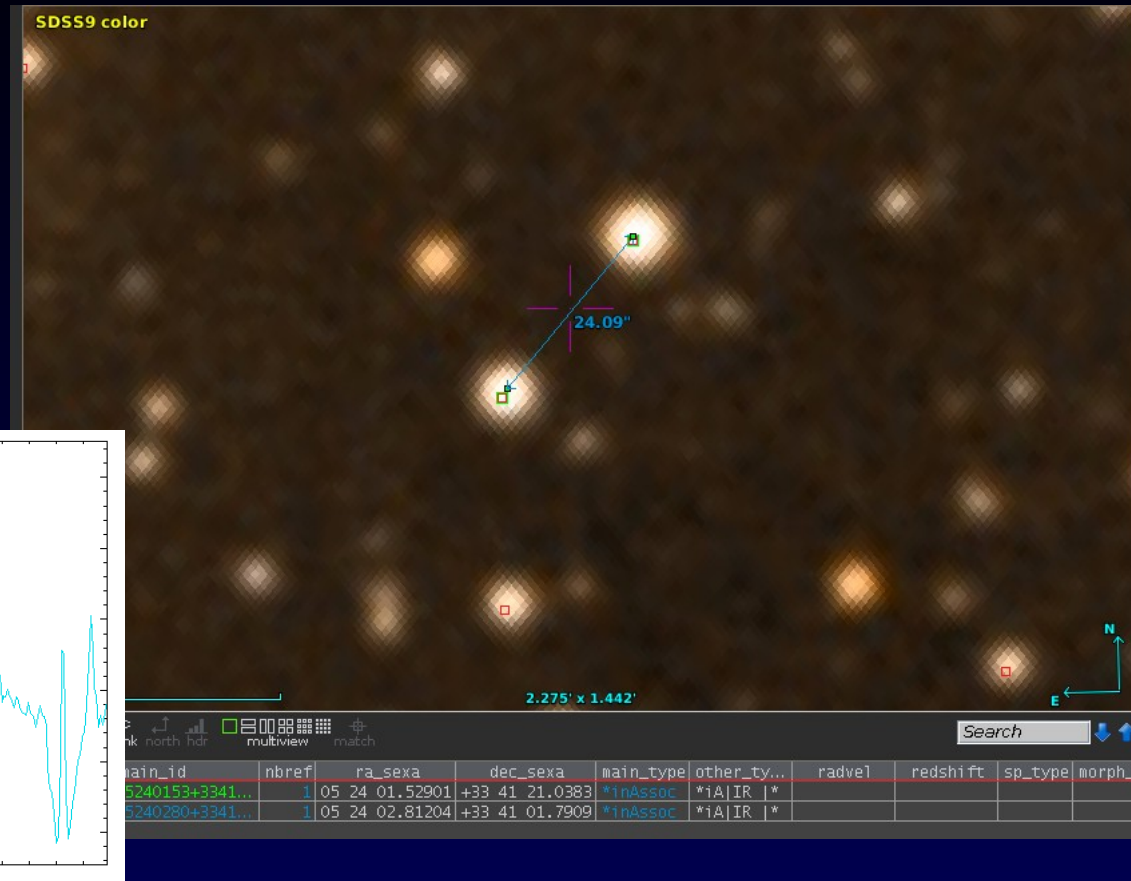
2 stars that are 24 arcsec apart

LAMOST J052402.81+334101.7

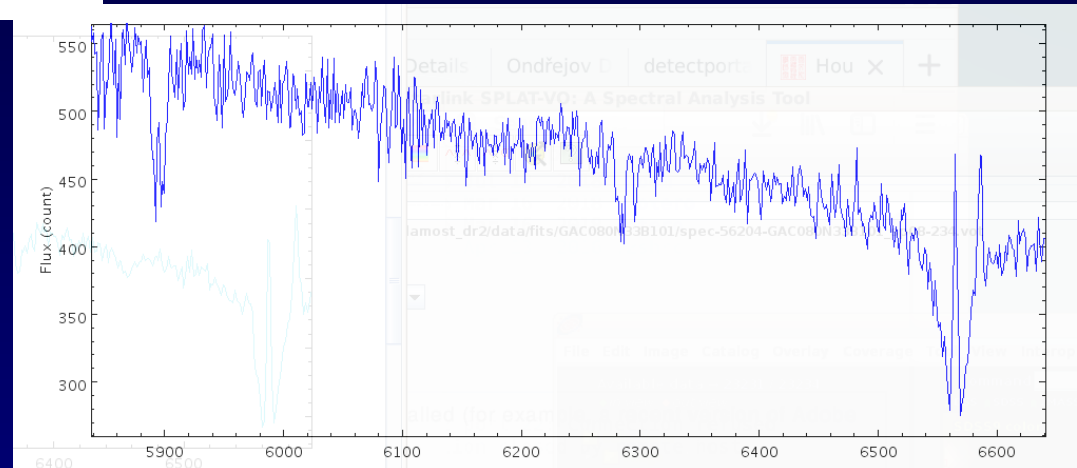
LAMOST J052401.53+334120.9

2MASS J05240280+3341017

2MASS J05240153+3341210



Is it the fiber light leakage ???  
OR Nebula lines ??



LAMOST J034912.80+240820.0  
Is Pleione 22 arcsec apart

# Confusion in Unique Identification

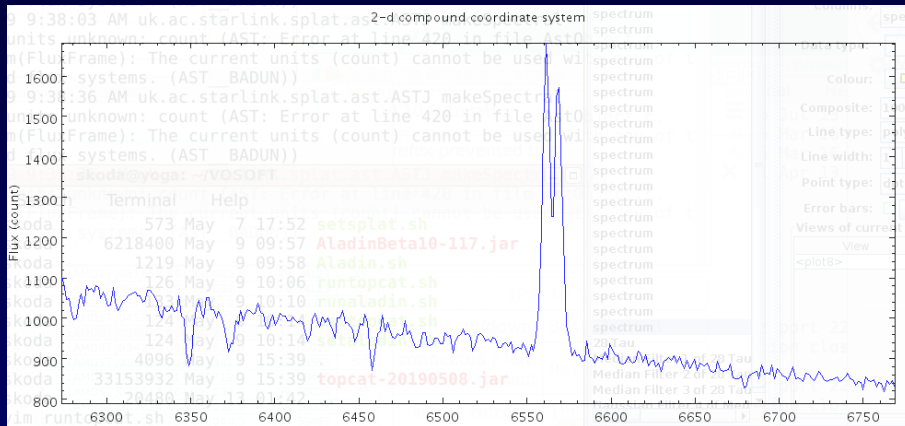
LAMOST J034912.80+240820.0

is Pleione (5mag) 22 arcsec apart

## Basic data :

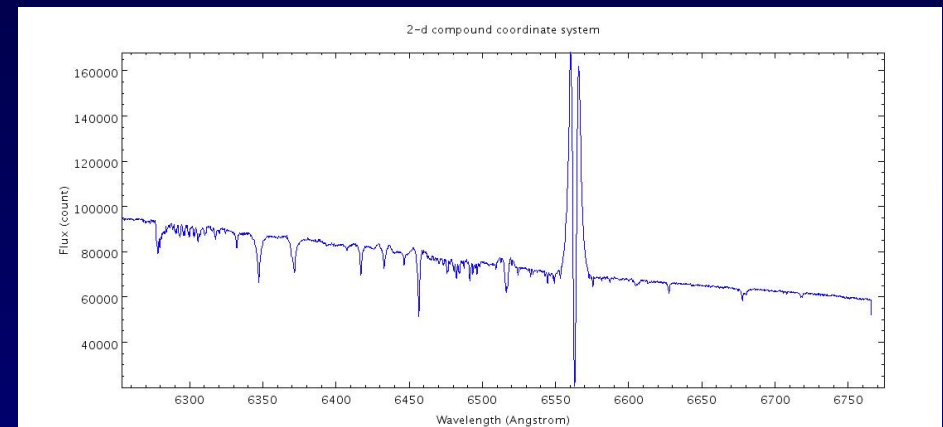
**LAMOST J034912.80+240820.0** -- Peculiar Star

Other object types: **Pe\*** (Ref)  
**ICRS coord. (ep=J2000) :** 03 49 12.800 +24 08 20.04 (Optical) [ ] D 2015MNRAS.449.1401H  
**FK4 coord. (ep=B1950 eq=1950) :** 03 46 14.053 +23 59 13.23 [ ]  
**Gal coord. (ep=J2000) :** 166.959918 -23.163713 [ ]  
 Spectral type: **A1mF1 D 2015MNRAS.449.1401H**



LAMOST MJD 56295

Ondrejov at MJD 56153





# Joining ML and VO

- Machine learning – supervised (labels)
  - Classical approach – preprocessing , table (unique ID)
    - Interpretation ID to original data
  - Requires enough labels - in science difficult
  - Deep learning – requires A LOT of data/labels
  - Active learning - needs to decide ON THE FLY
  - RECLASSIFICATION – interactive
  - Oraculum – needs COMPLEX information to decide
  - METADATA, comparable DATA
- VO Interoperability – crucial part of process**

# Joining ML and VO

- ML is Big Data = Done in cloud – GPU cluster, Python
- Interactive visualization + classification GUI

VO-CLOUD DETAILS OF JOB

Home Manage filesystem Jobs Download history Create job Jupyter Settings Admin Help Logout (skoda)

**active-learning-test(copy)(copy)**

Type	Id	Phase	Worker	Created	Started	Finished
Active_learning	156-357	COMPLETED	local worker	4/24/19 7:47:17 AM	4/24/19 7:47:17 AM	4/24/19 7:47:19 AM

Run again Delete

Preview

index.html - Fullscreen

## Spectra

Name	Ra	Dec	Mag	Prediction	Label	Iteration
	NEW	3.82002777549371	-	double peak	not sure	1

1-single peak  
 2-double peak  noteworthy    
 3-not sure  
 4-bad

6530 6540 6550 6560 6570 6580 6590 6600 6610

# Conclusions

- Active learning overcomes the lack of labeled data
- ML on big spectra archives may identify new interesting objects yet **unknown/unexpected**
- Domain adaptation gives LABELS on unlabeled set
- Crucial is interactive visualization of candidates .....
- Deep learning shows its strength
  
- Active Learning - ORACULUM requires **VO** to have  
**METADATA** to decide correctly  
**OTHER DATA** (global interoperability !)
- ML needs to **visualize** data as part of its process now!

**THANK YOU**