

# Deep-learning for galaxy morphology and evolution

**A review of on-going projects - questions / ideas**

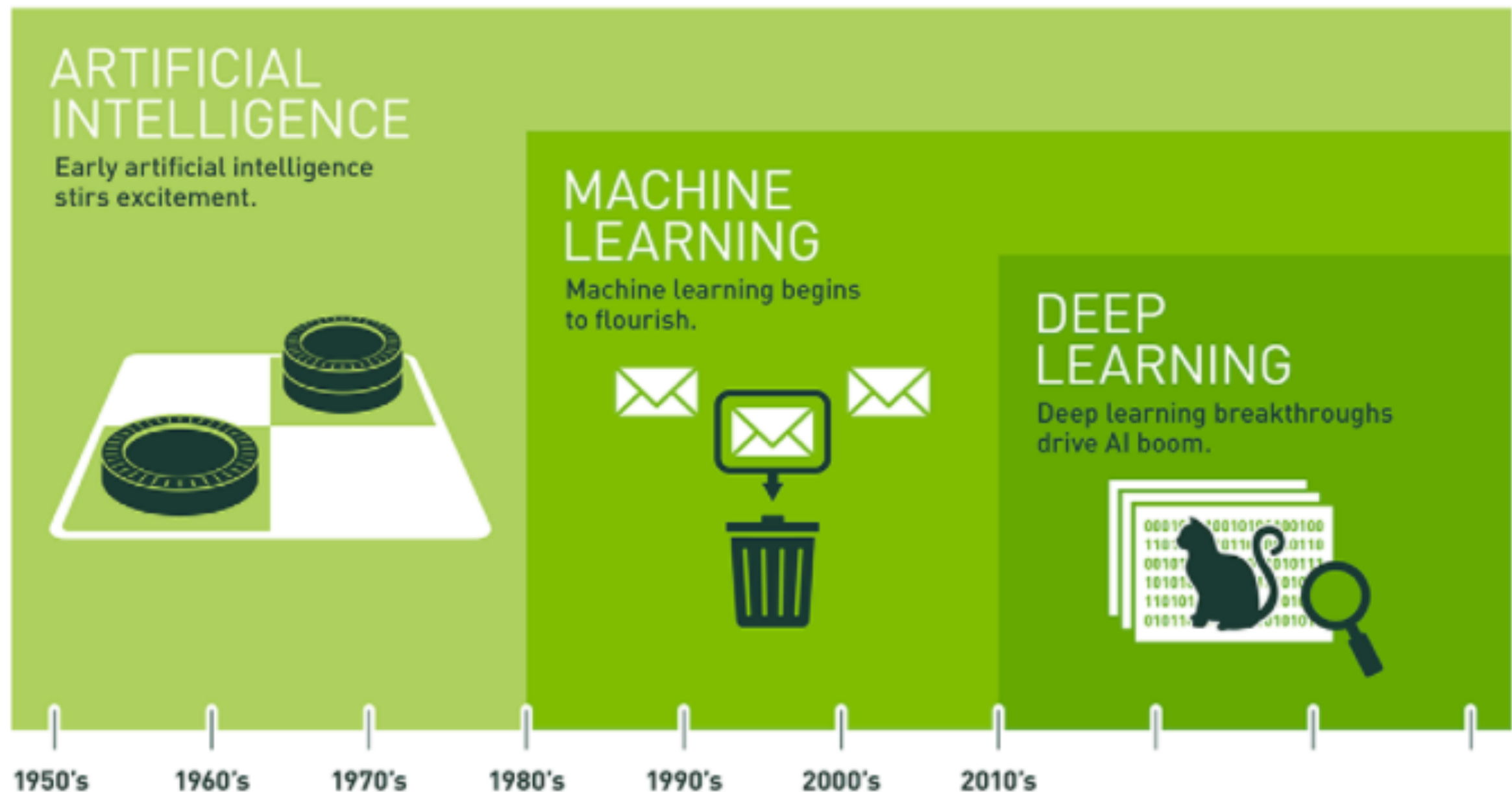
Marc Huertas-Company

D. Koo, J. Primack, H. Dominguez-Sanchez, M. Bernardi, S. Faber, F. Caro,  
D. Tuccillo, C. Lee, B. Margalef-Bentabol, E. Decencière, S. Velasco-Forero, G. Cabrera-Vives....

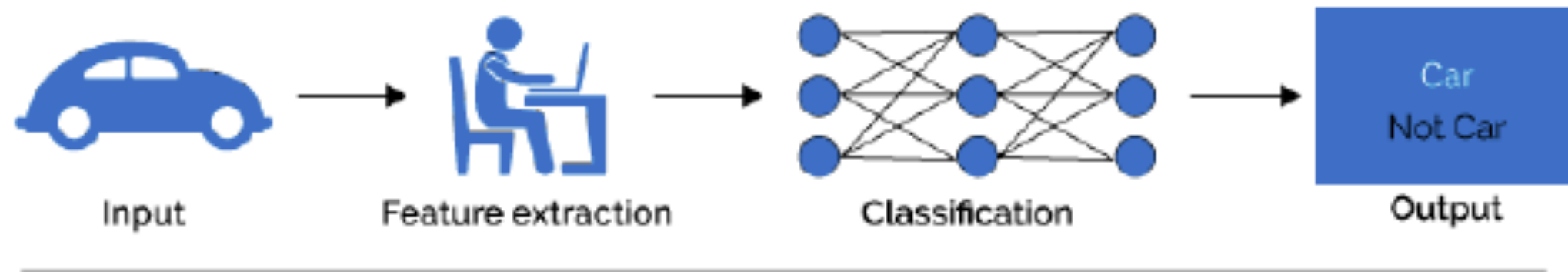
# **GOLDEN AGE FOR ARTIFICIAL INTELLIGENCE** with the generalization of deep-learning since 2012

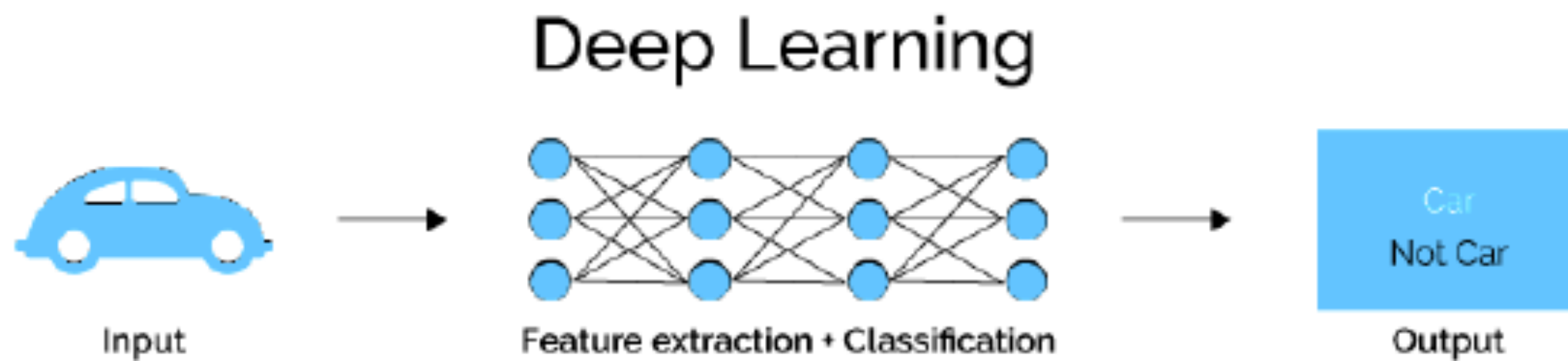
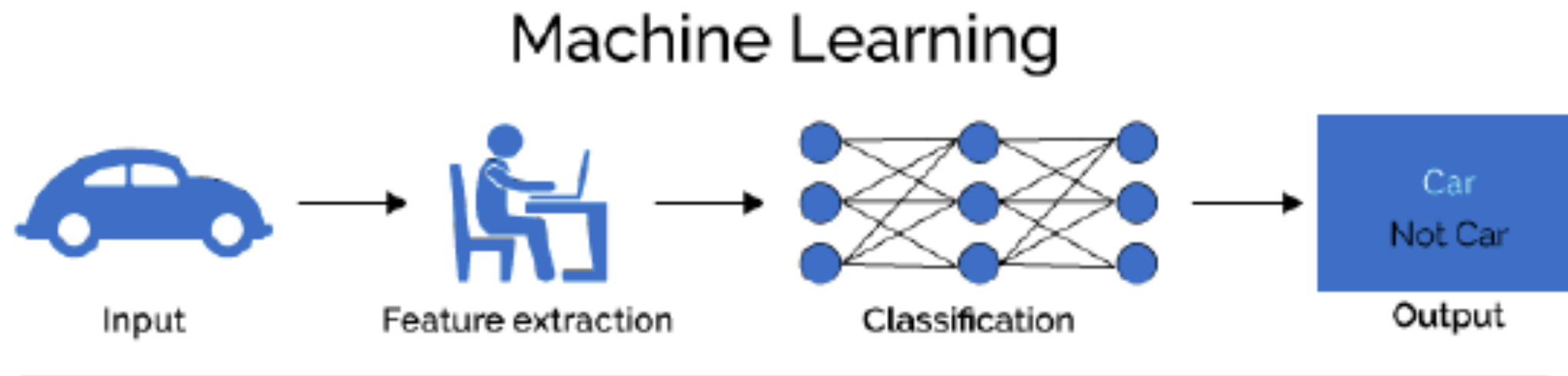
**(Nature, 01/2016)**

**“Deep learning is killing every problem in AI.”**



# Machine Learning





The 3 “phases” of deep-learning:

**skepticism, acceptance, frustration**

# **GOLDEN AGE FOR ARTIFICIAL INTELLIGENCE**

(“Artificial Intelligence”) with the generalization of deep-learning since 2012

## **QUESTIONS WE ARE FOCUSING ON :**

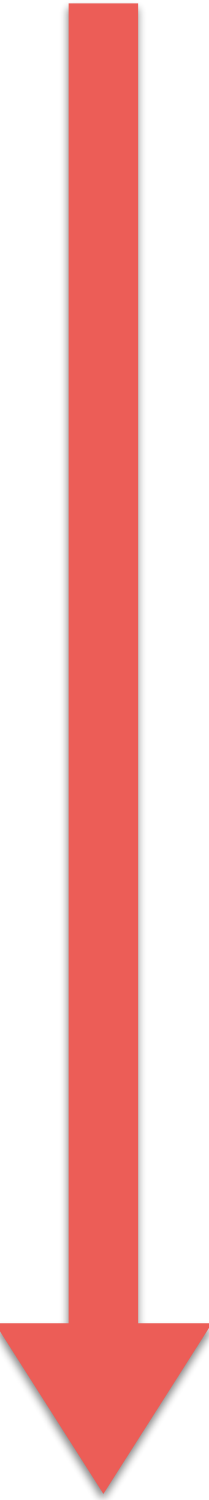
How AI (i.e. deep-learning) can be used to understand galaxy formation?

Can we do things with AI that we could not do before?

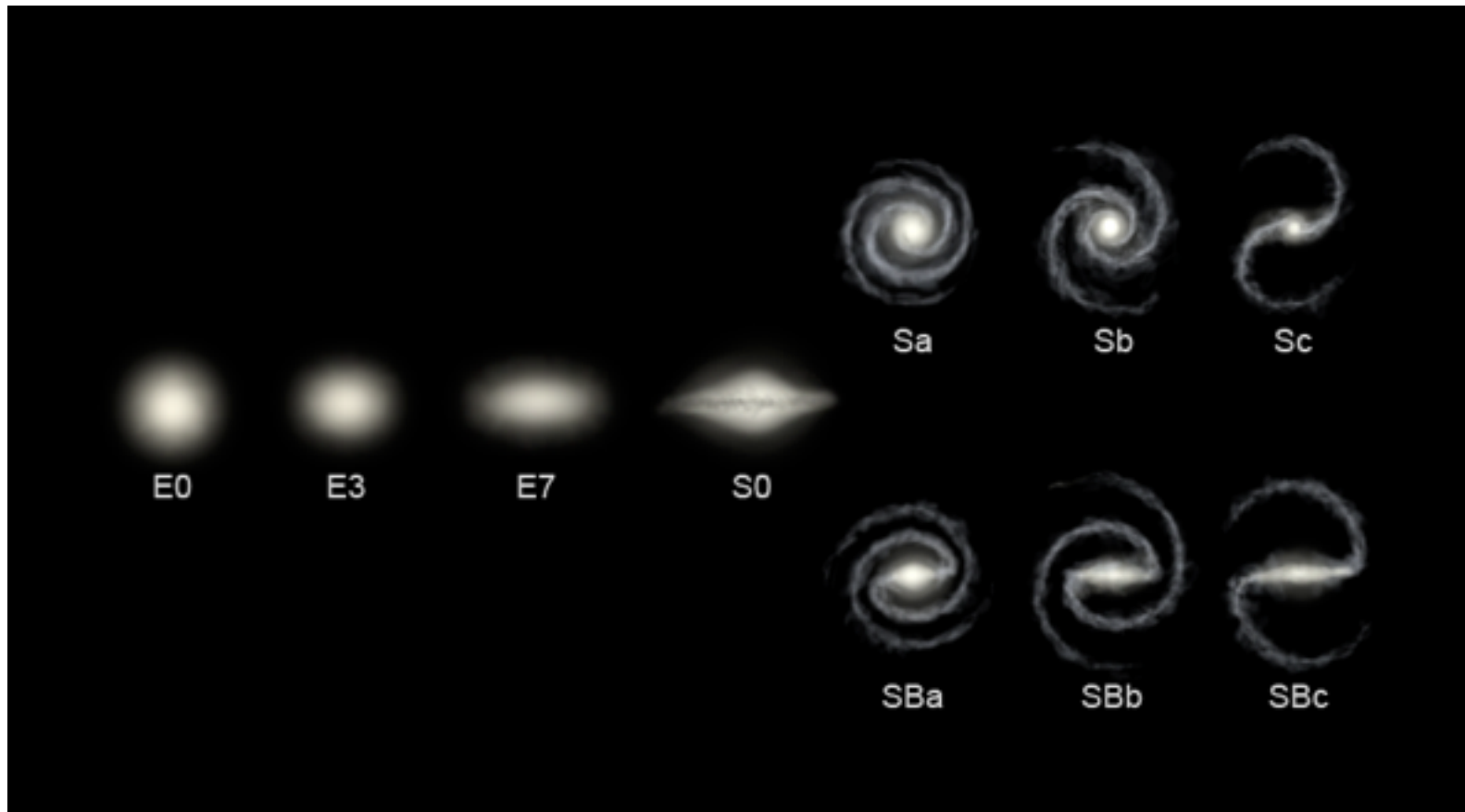
Can we learn something new about the physics of galaxies?

# DL FOR GALAXIES?

- **GROUP #1:** Time consuming tasks that humans do easily but classically challenging for computers - classification of objects
- **GROUP #2:** Efficient and fast quantitative measurements on large amount of (multi-lambda) data [photoz's, sizes, ellipticities]
- **GROUP #3:** Find hidden new observables in the data, - Linking observations and theory
- **GROUP #4:** Finding the unknown?



- **GROUP #1:** Time consuming tasks that humans do easily but classically challenging for computers - classification of objects

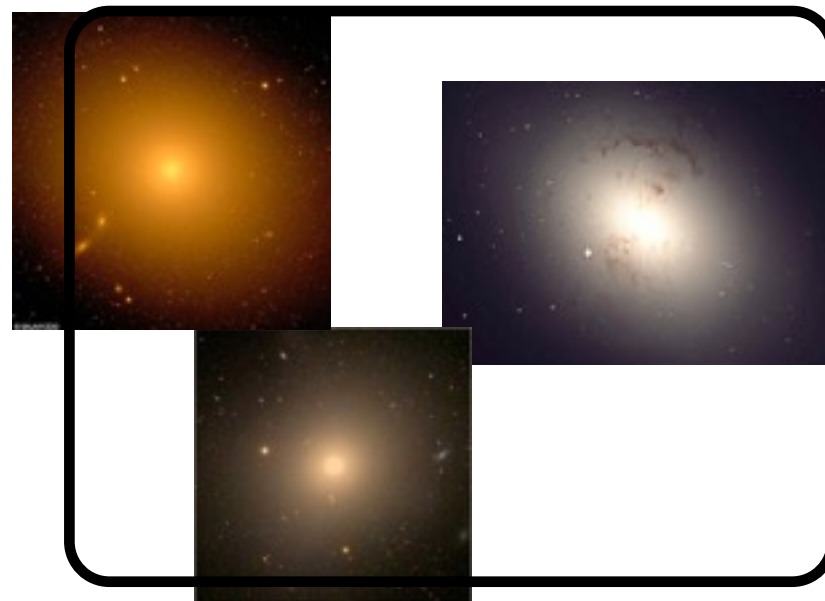


# The Hubble Sequence

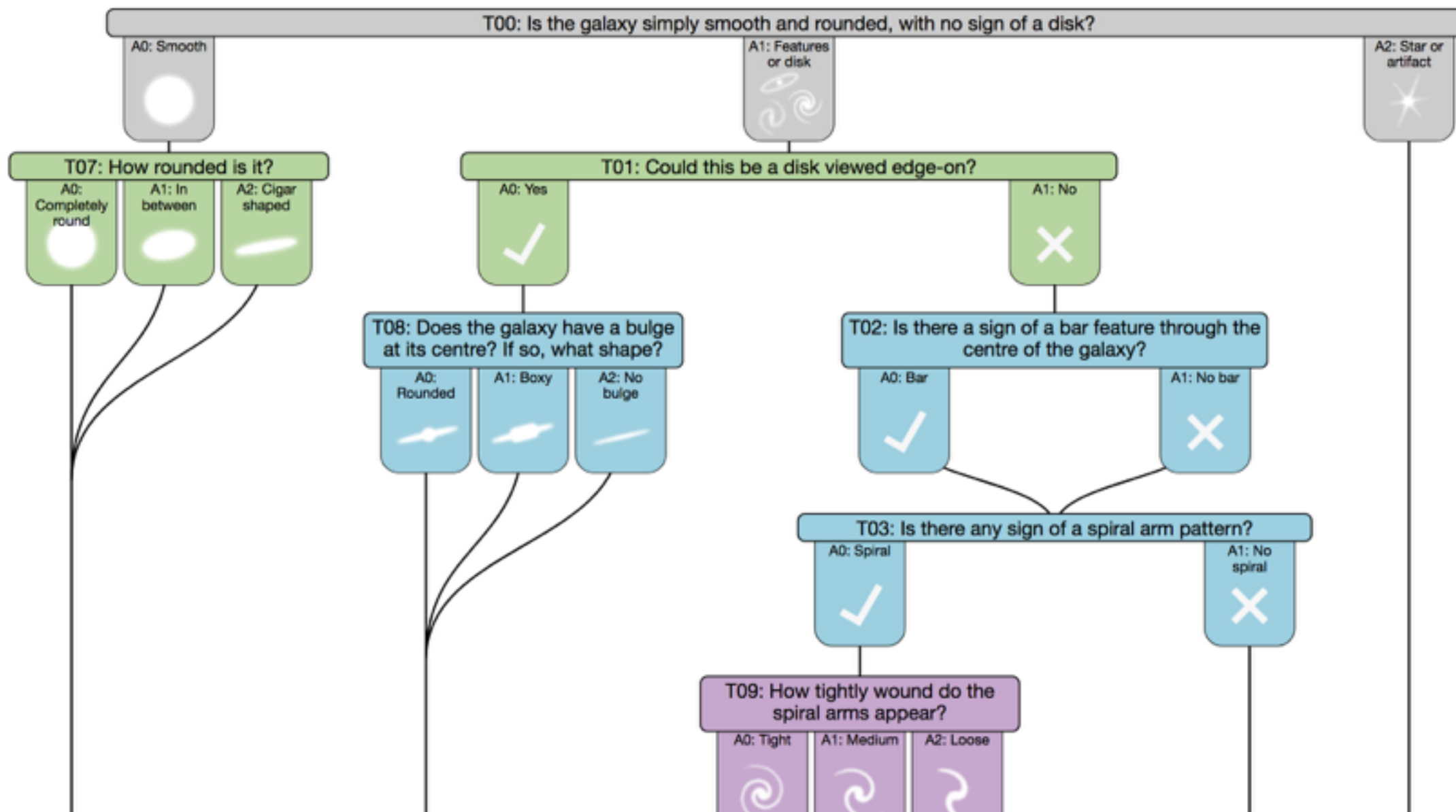
**BOX I**

**BOX II**

“Objects in the same box experienced the same physics”



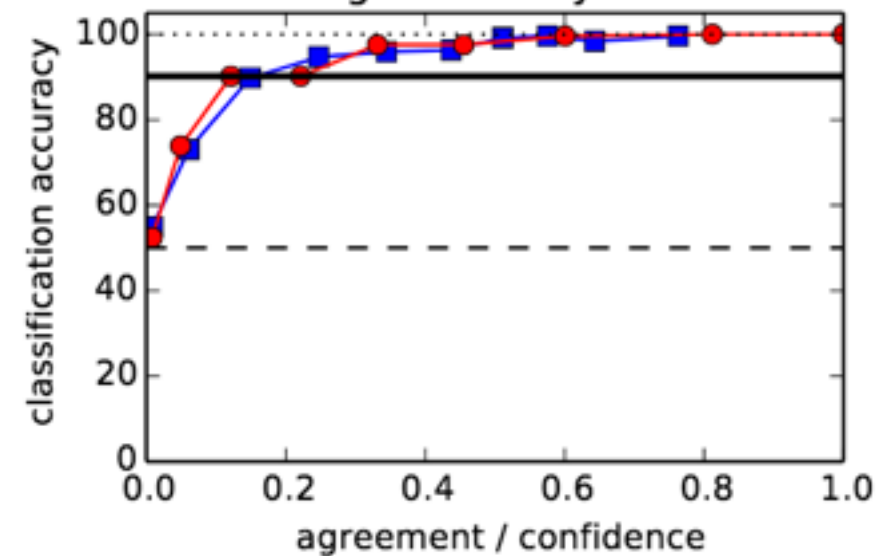
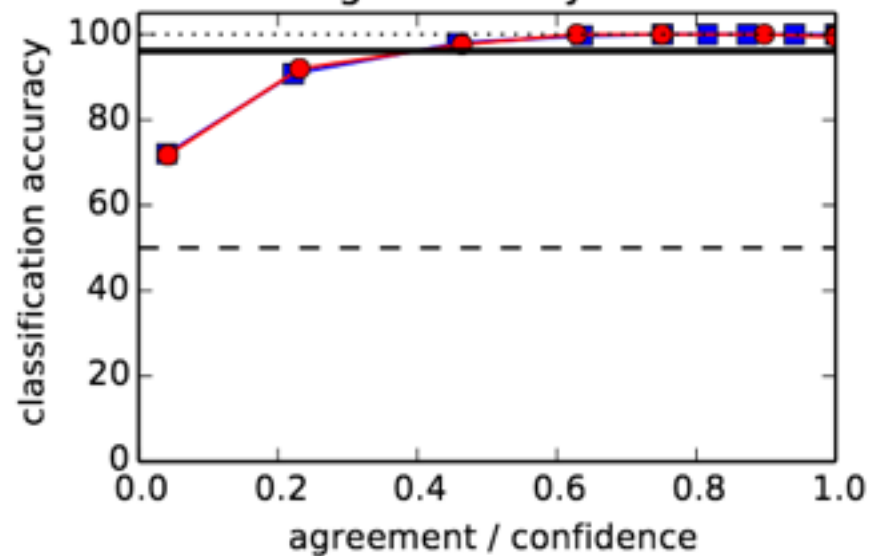
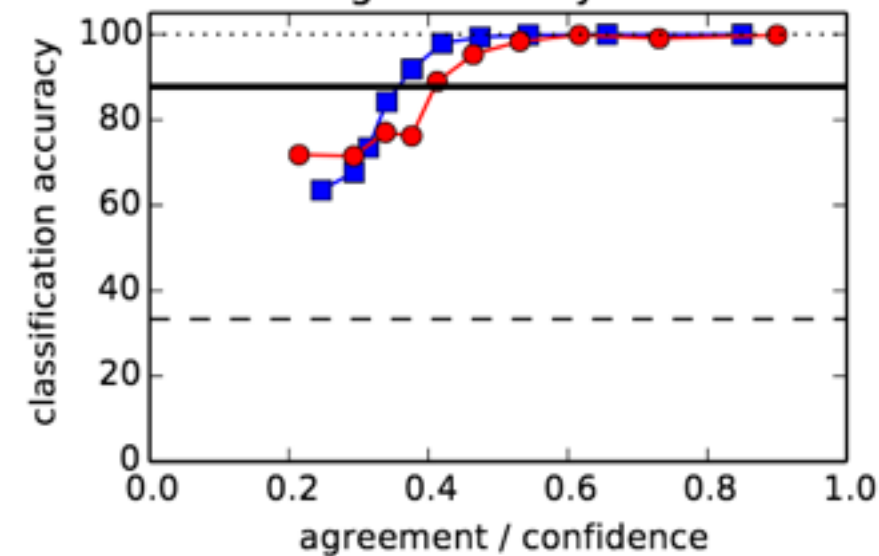




Q1: smoothness, 6144 examples  
average accuracy: 87.79%

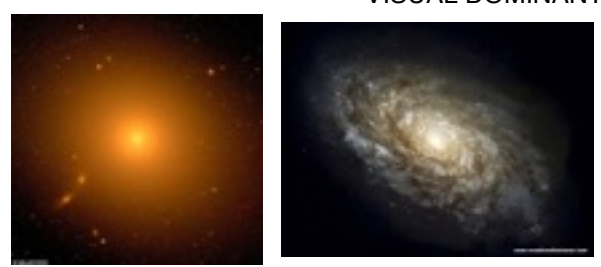
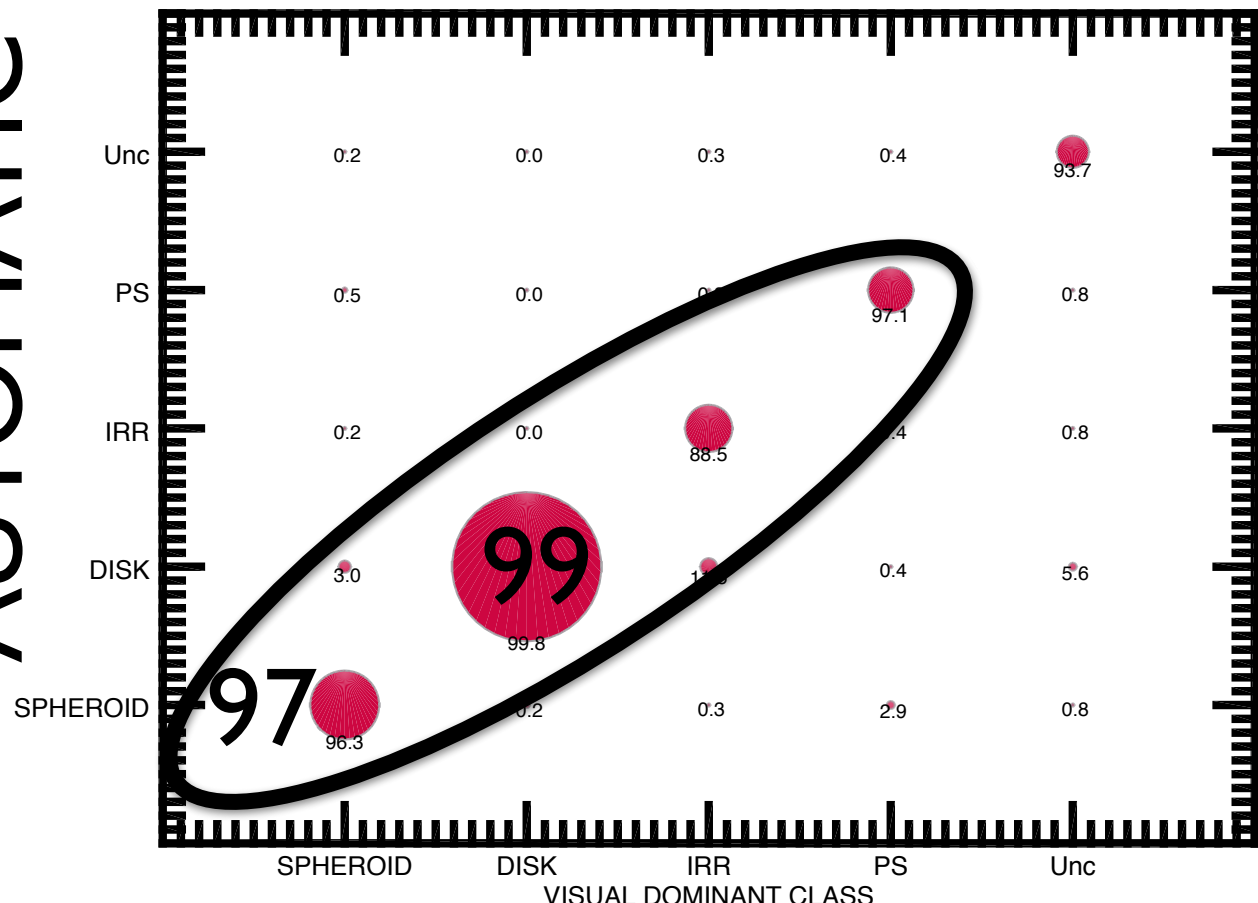
Q2: edge-on, 3362 examples  
average accuracy: 96.04%

Q3: bar, 2449 examples  
average accuracy: 90.16%



# HIGH-REDSHIFT MORPHOLOGIES (CANDELS)

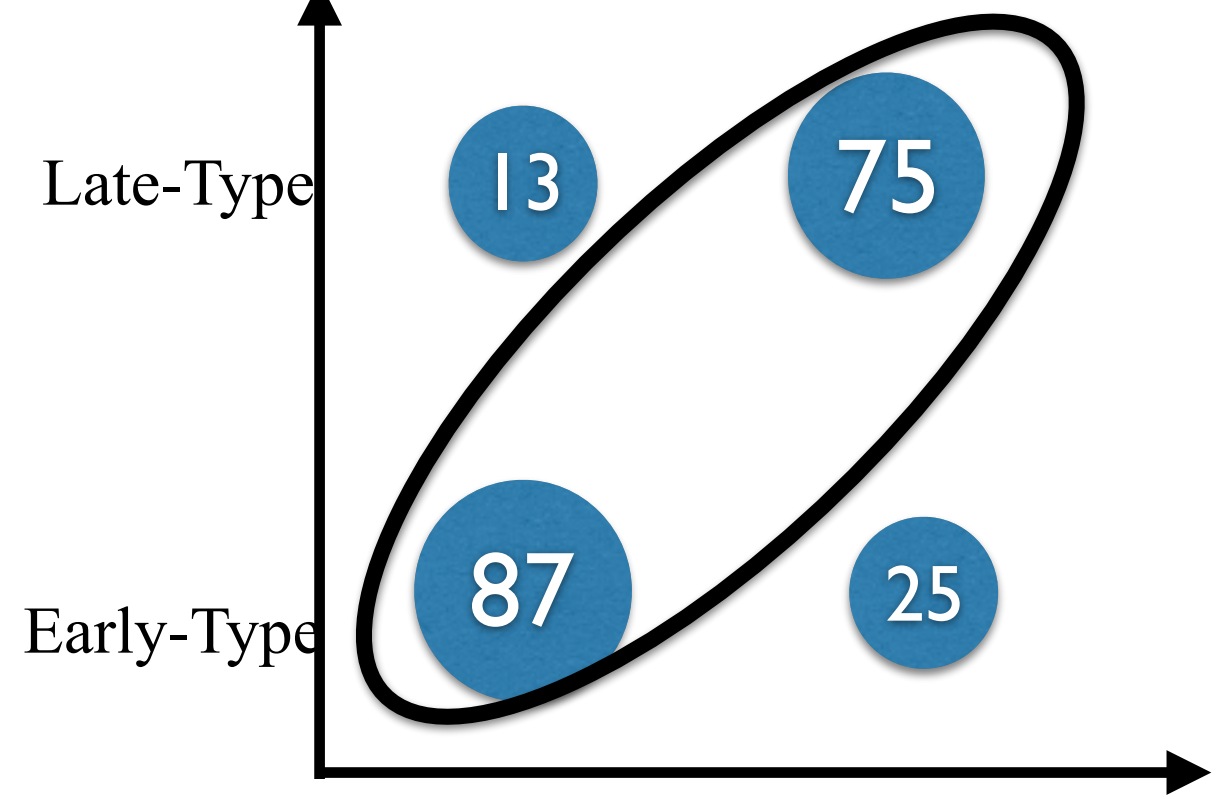
AUTOMATIC



**CNNs**

**VISUAL**

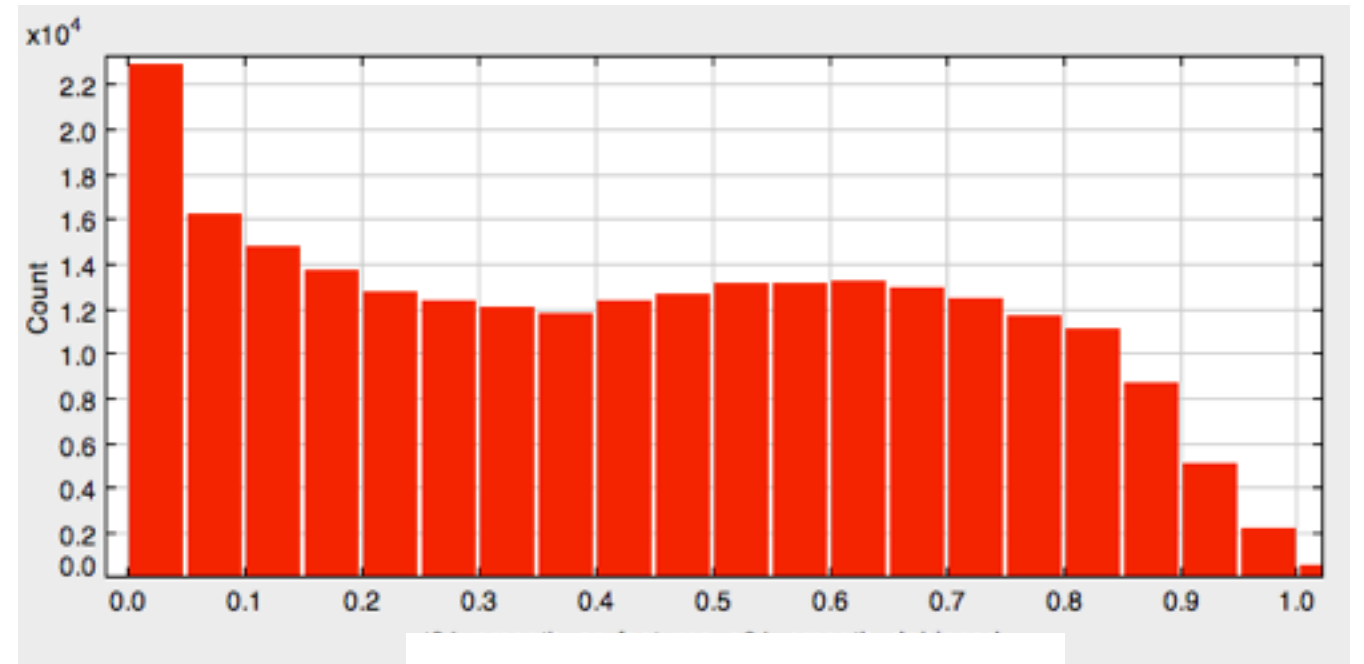
**AUTOMATIC**



**SVMs**

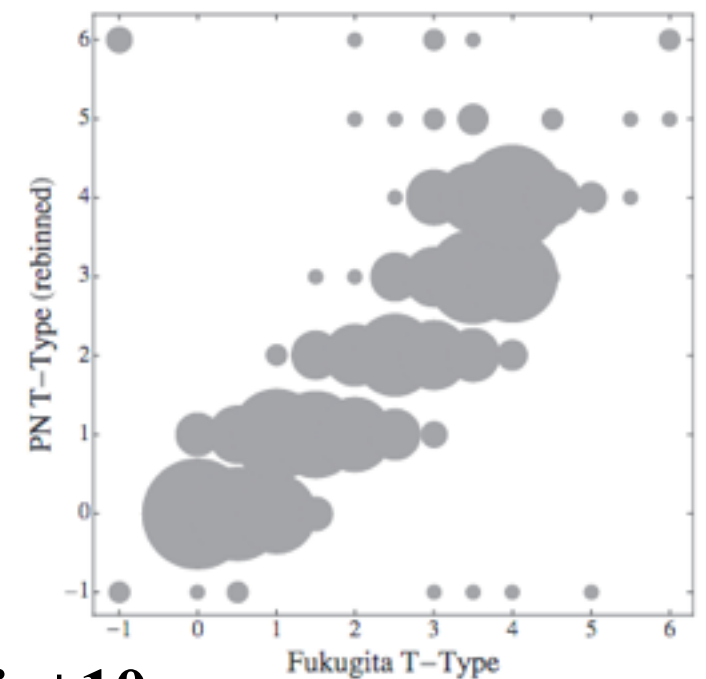
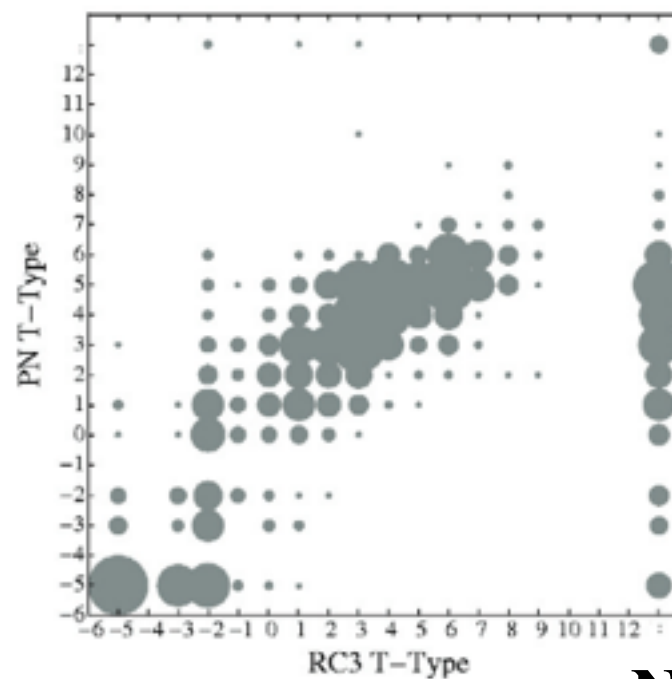
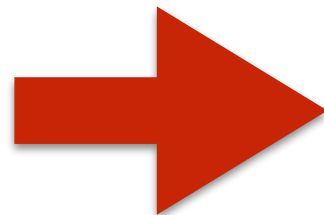
Can we improve human biases with machines?

GZOO



PROBABILITY

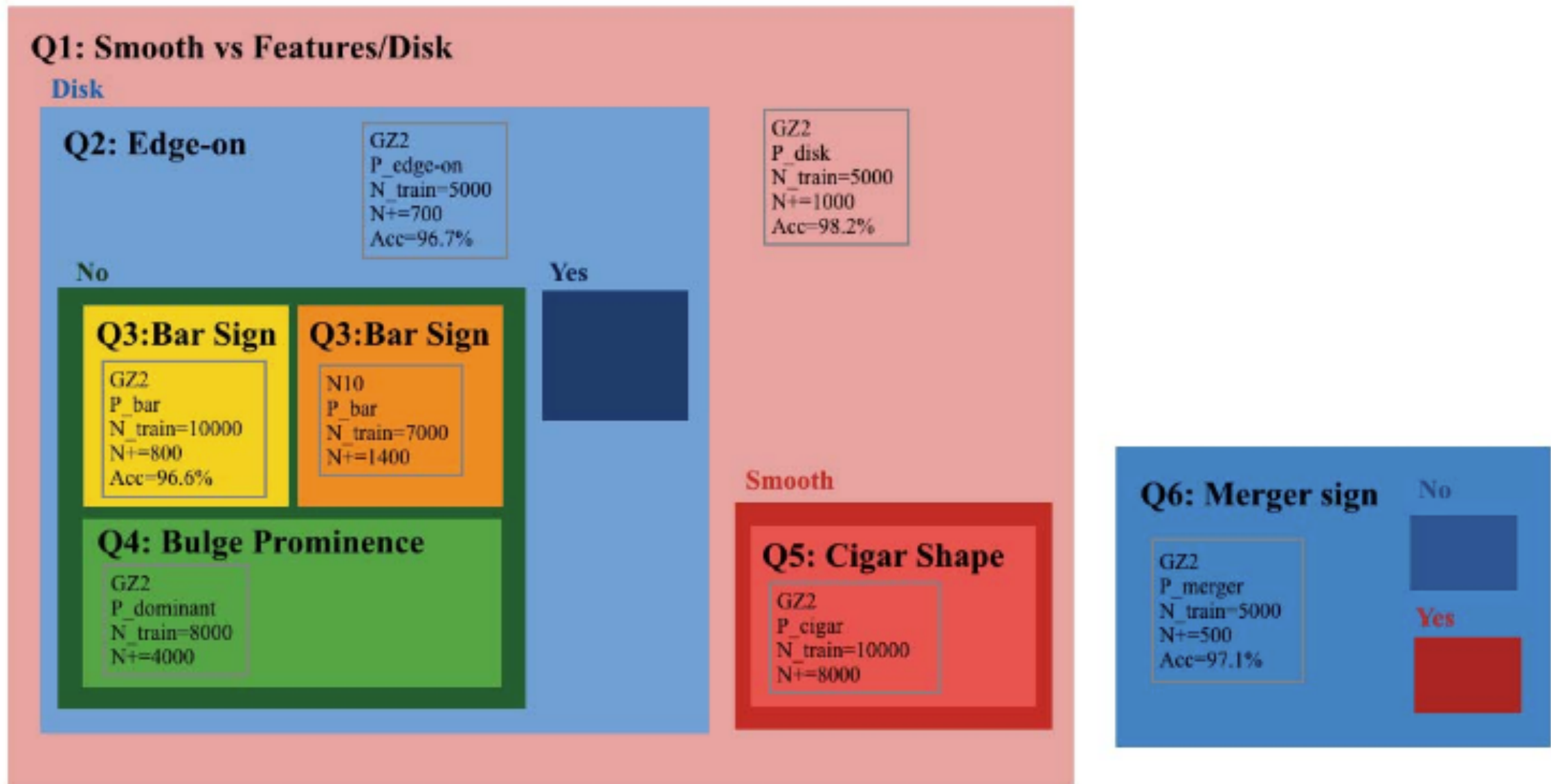
T-Type  
(astronomers)



Nair+10

Fukugita T-Type

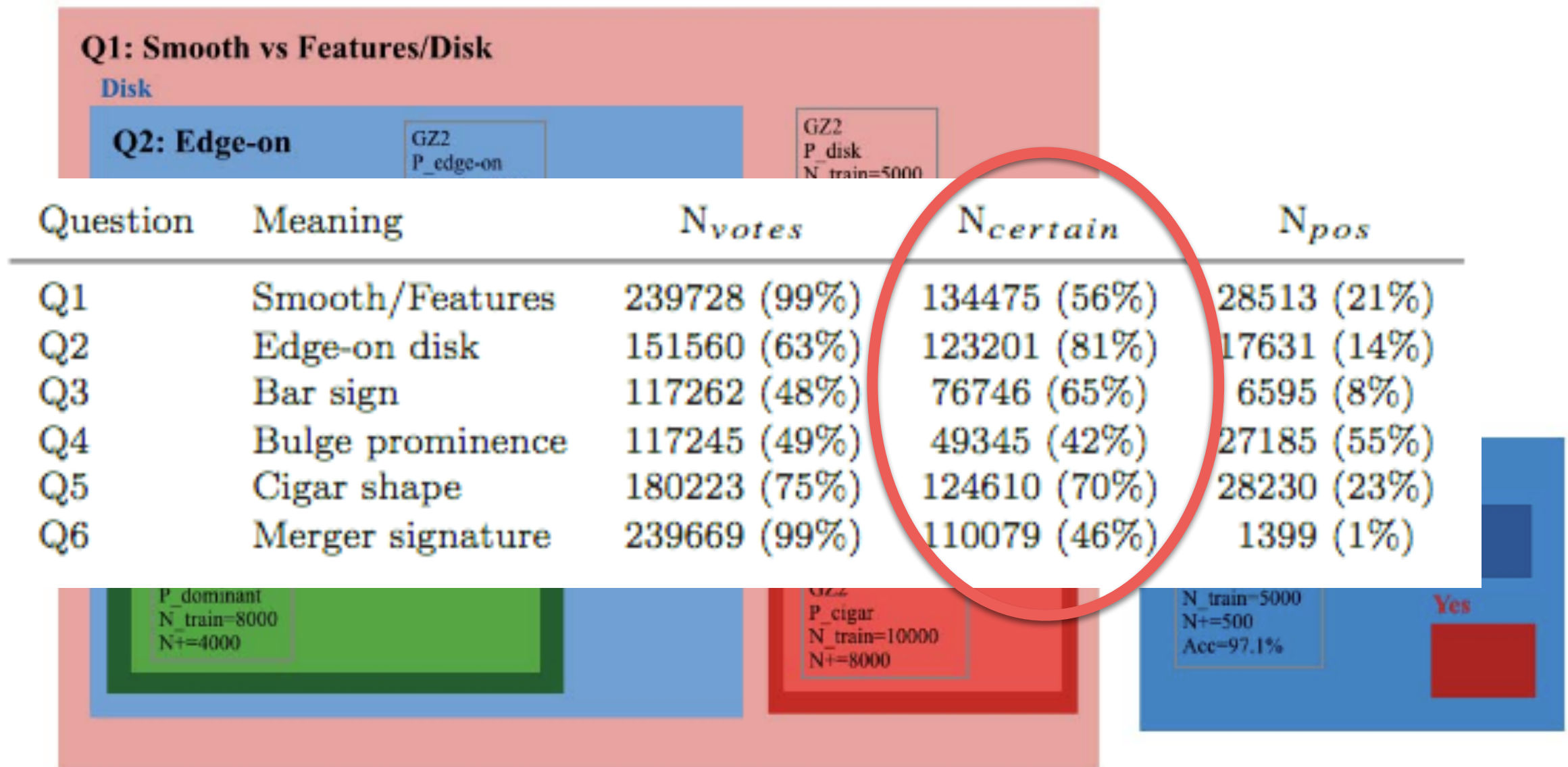
# REVISITING THE SDSS MORPHOLOGY



Select only “safe” classifications for training [ $N > 5$ ,  $P > 0.7$ ]  
Binary classification for each feature separately

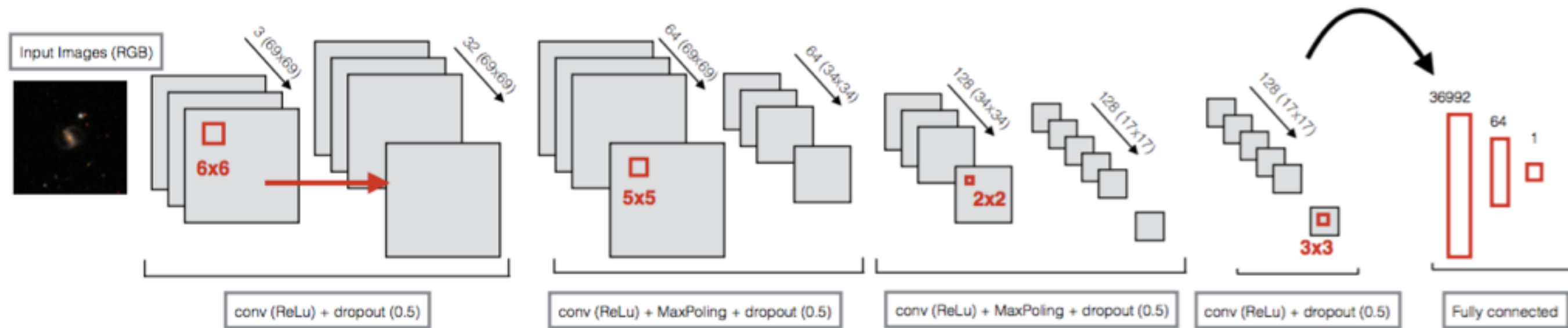


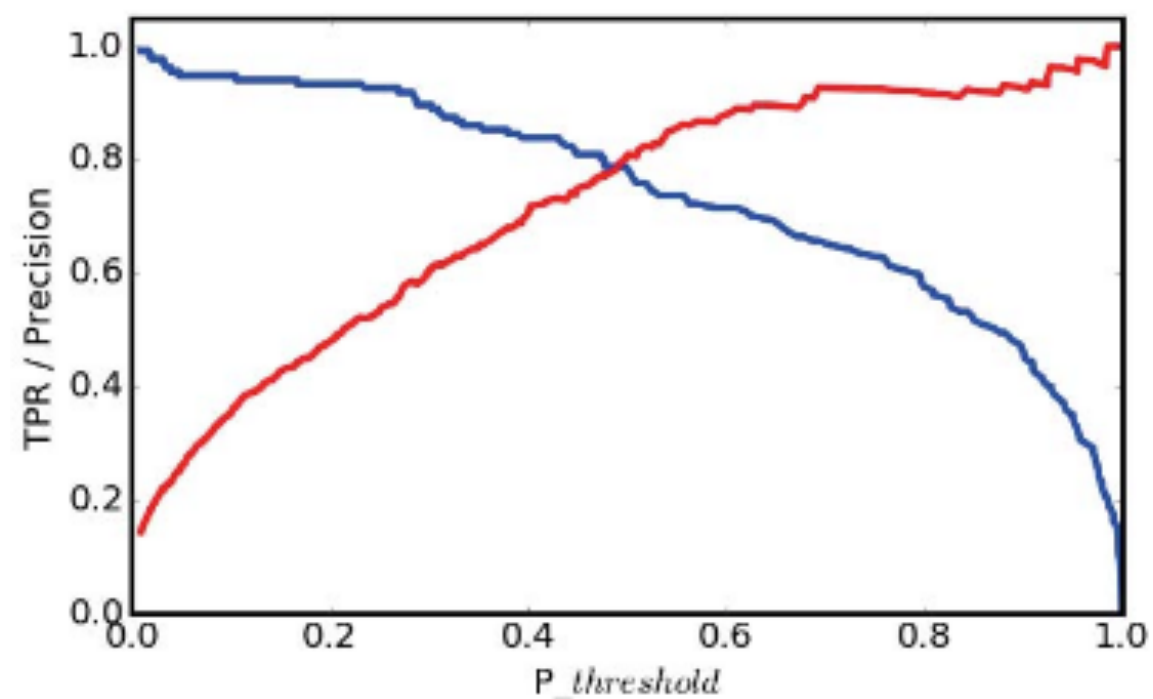
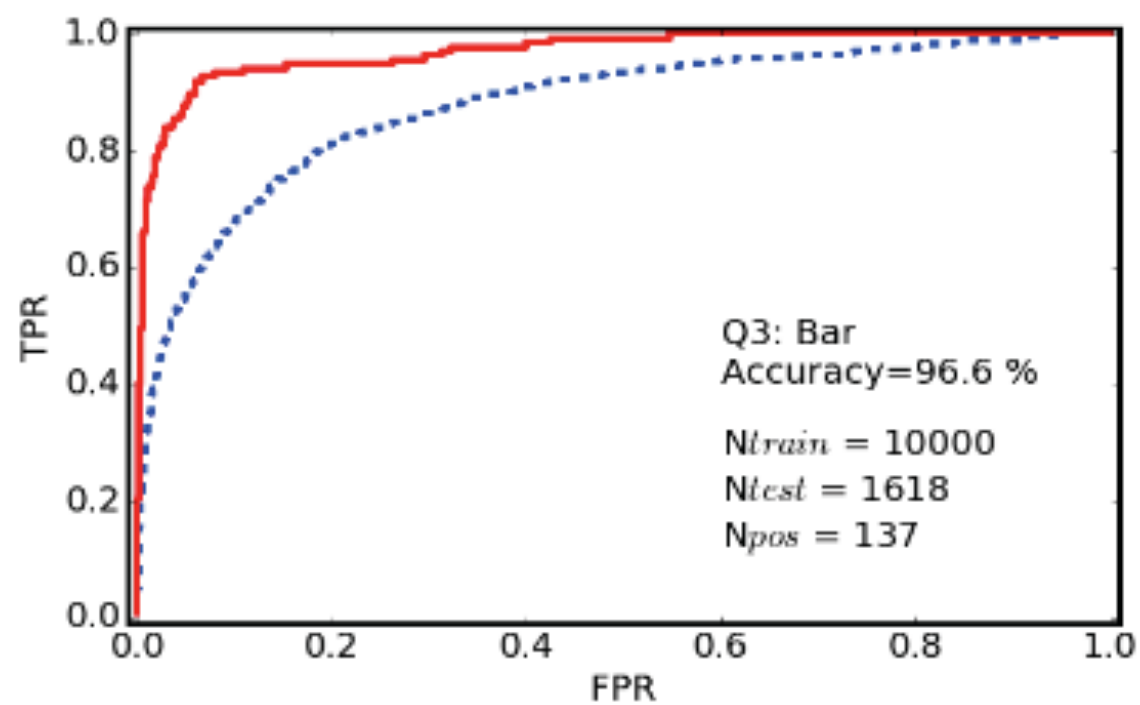
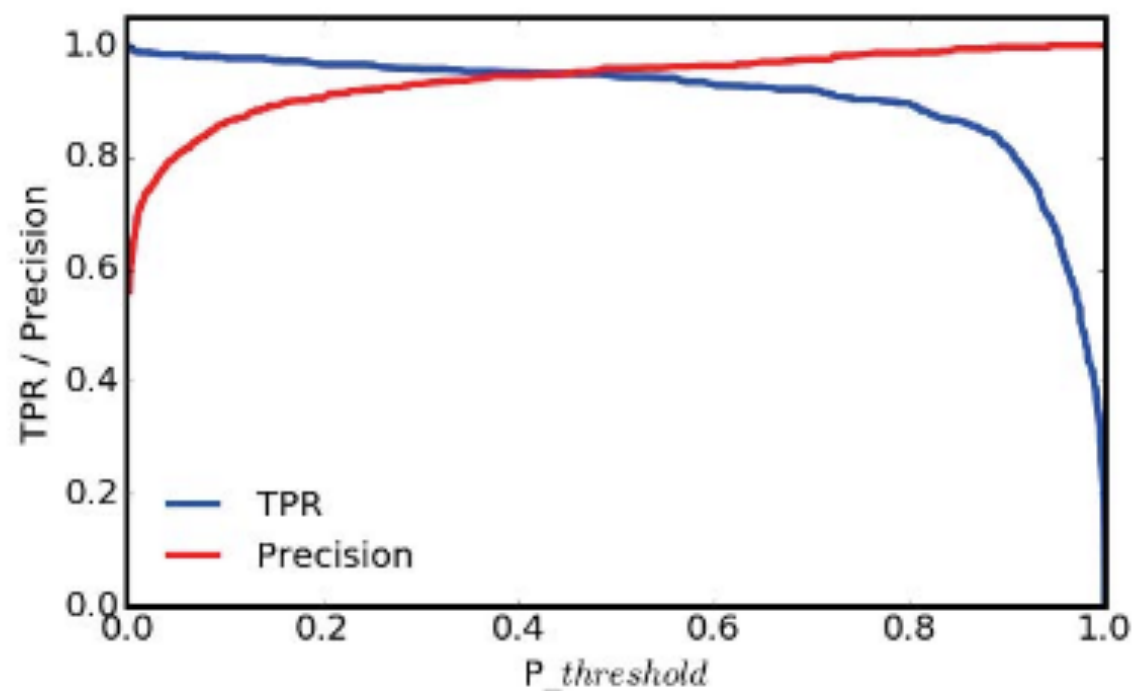
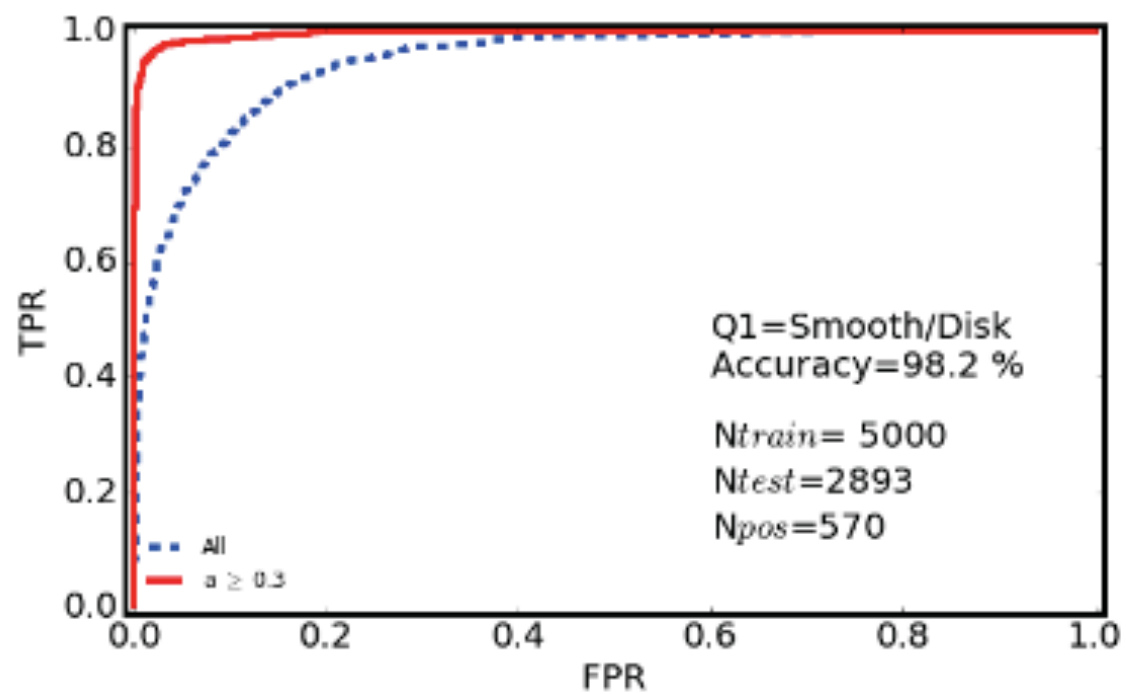
# REVISITING THE SDSS MORPHOLOGY



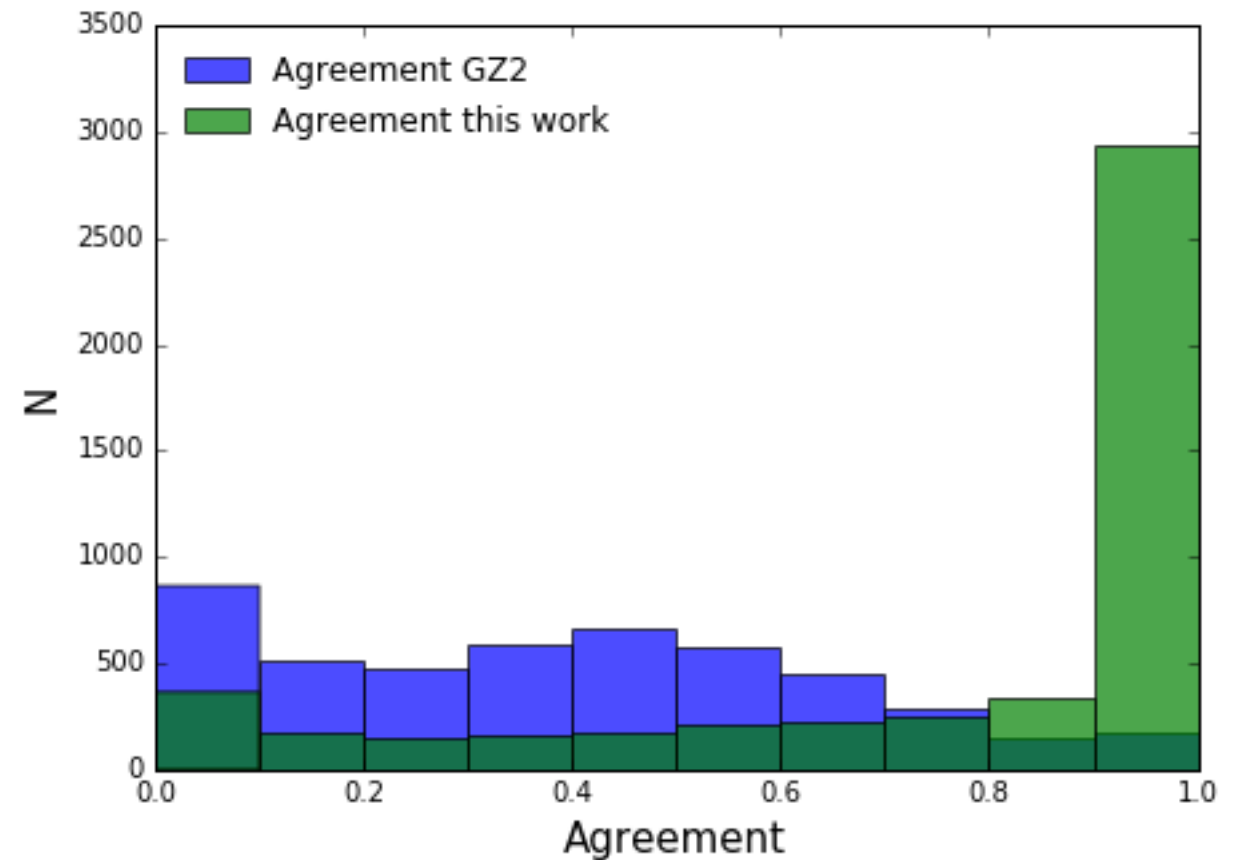
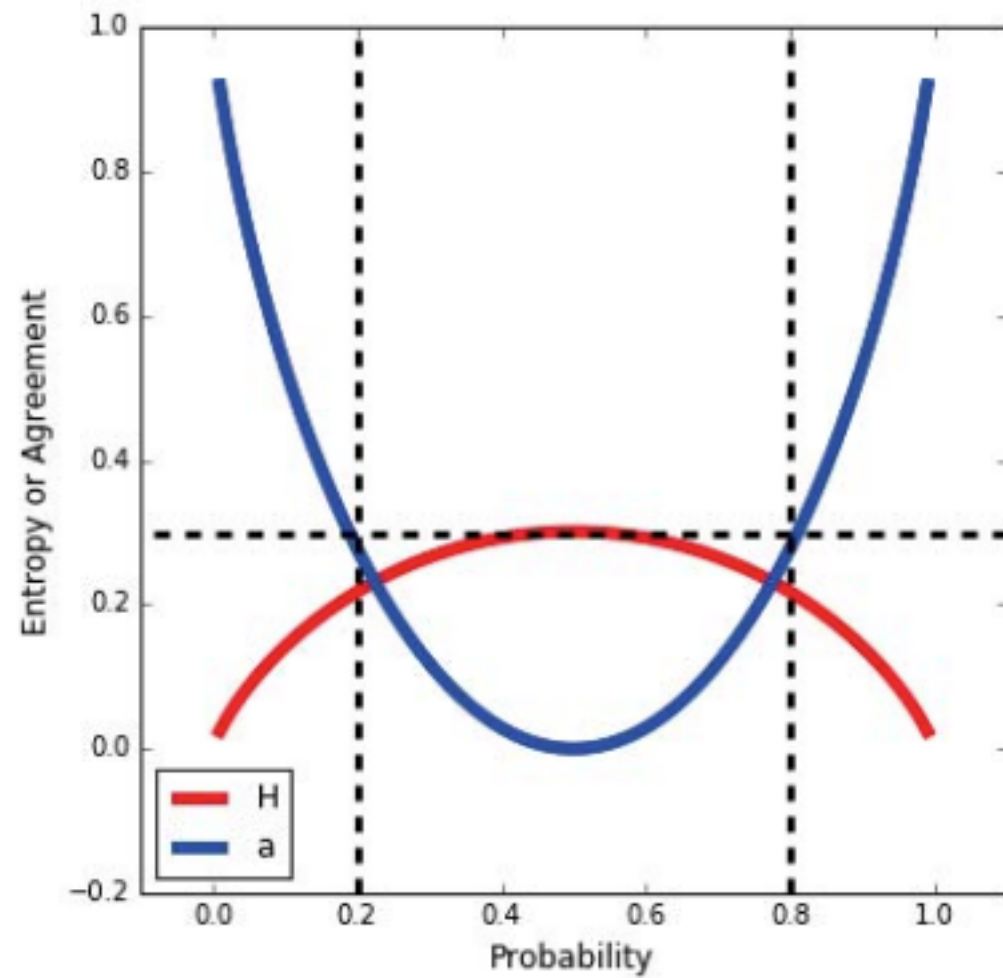
Select only “safe” classifications for training [ $N > 5$ ,  $P > 0.7$ ]  
Binary classification for each feature separately

# VERY SIMPLE ARCHITECTURE

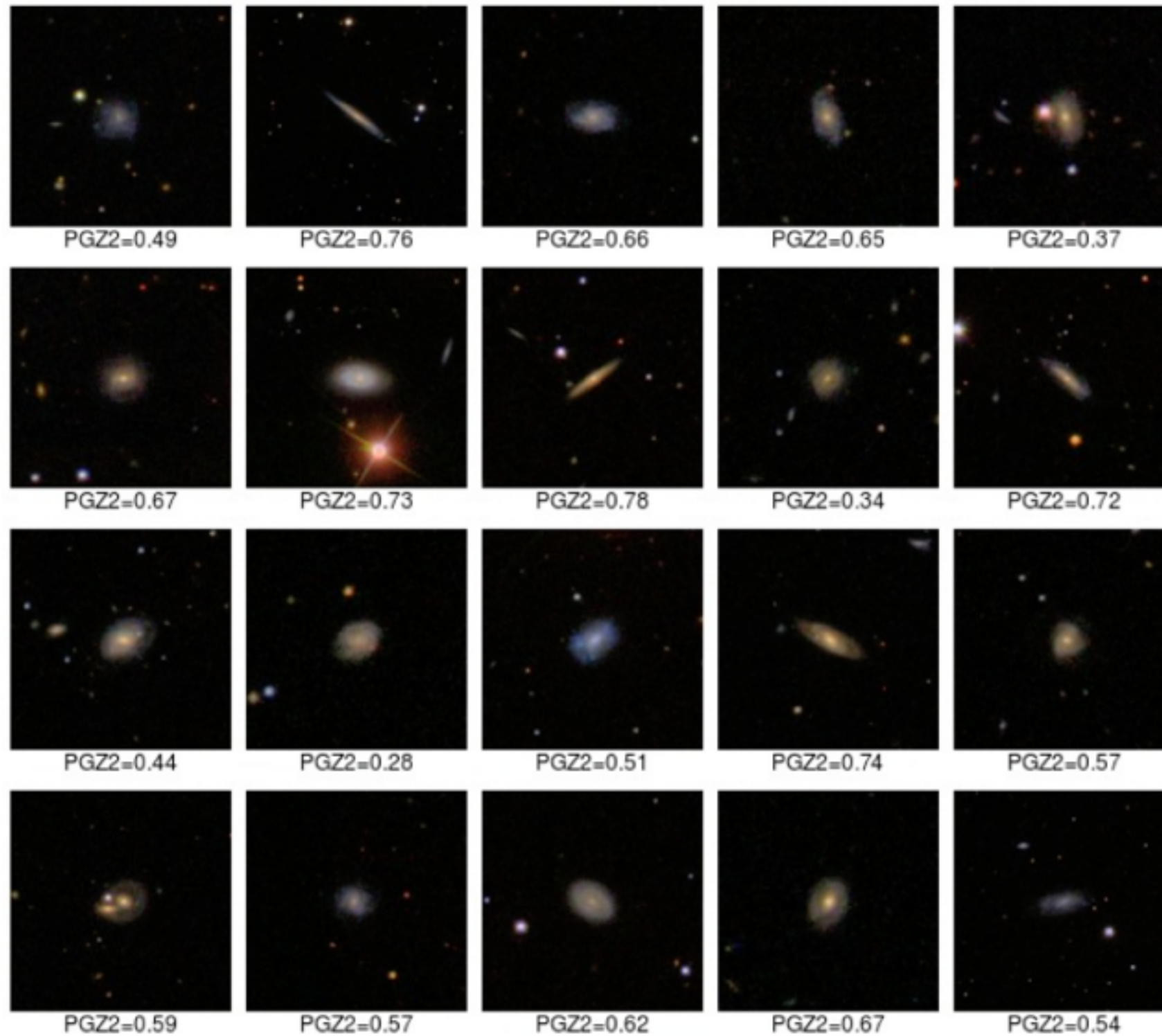




# Entropy of classification







**SECURE DISKS GALAXIES FOR DL  
- UNCLEAR FOR PEOPLE**

How robust to different datasets?  
Do we always need a big training  
set?

**DATA FROM  
NEW SURVEY**

How robust to different datasets?  
Do we always need a big training  
set?

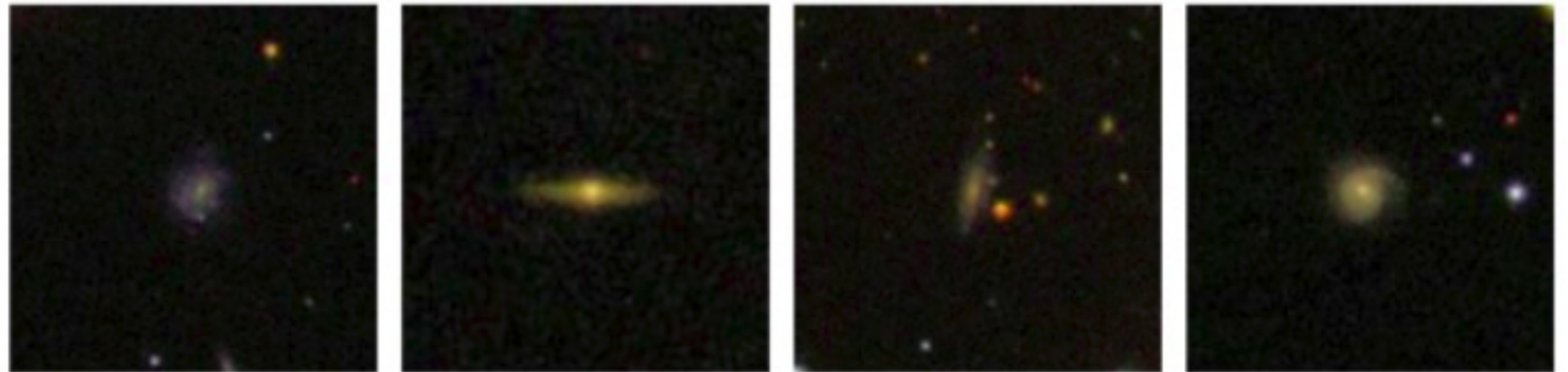
**Transfer knowledge?**

DEEP-LEARNING  
BASED  
MACHINE

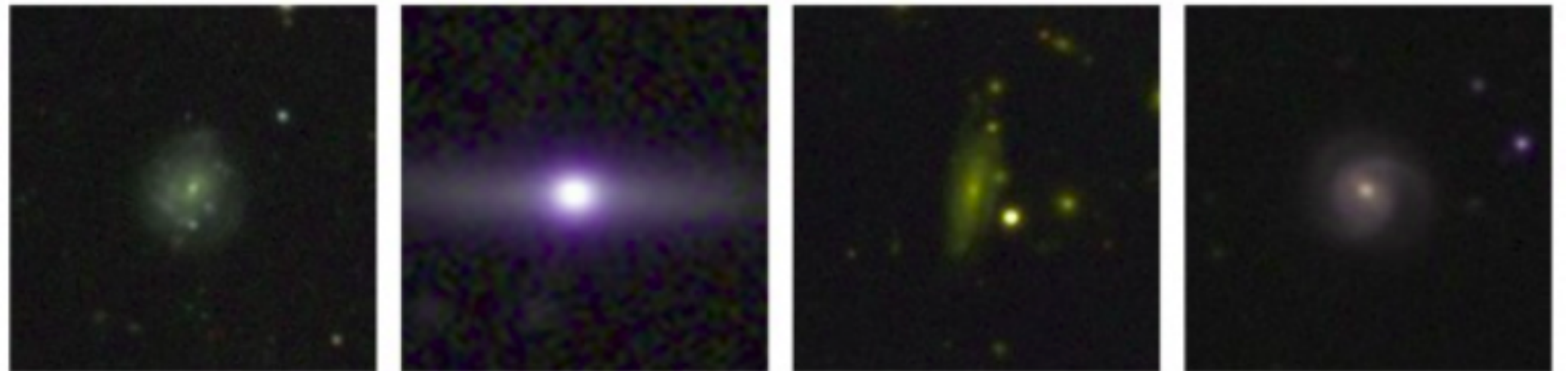
“Improved”  
Galaxy ZOO like  
classifications for  
for the entire  
sample

**Human classifications  
from existing survey**

**SDSS**

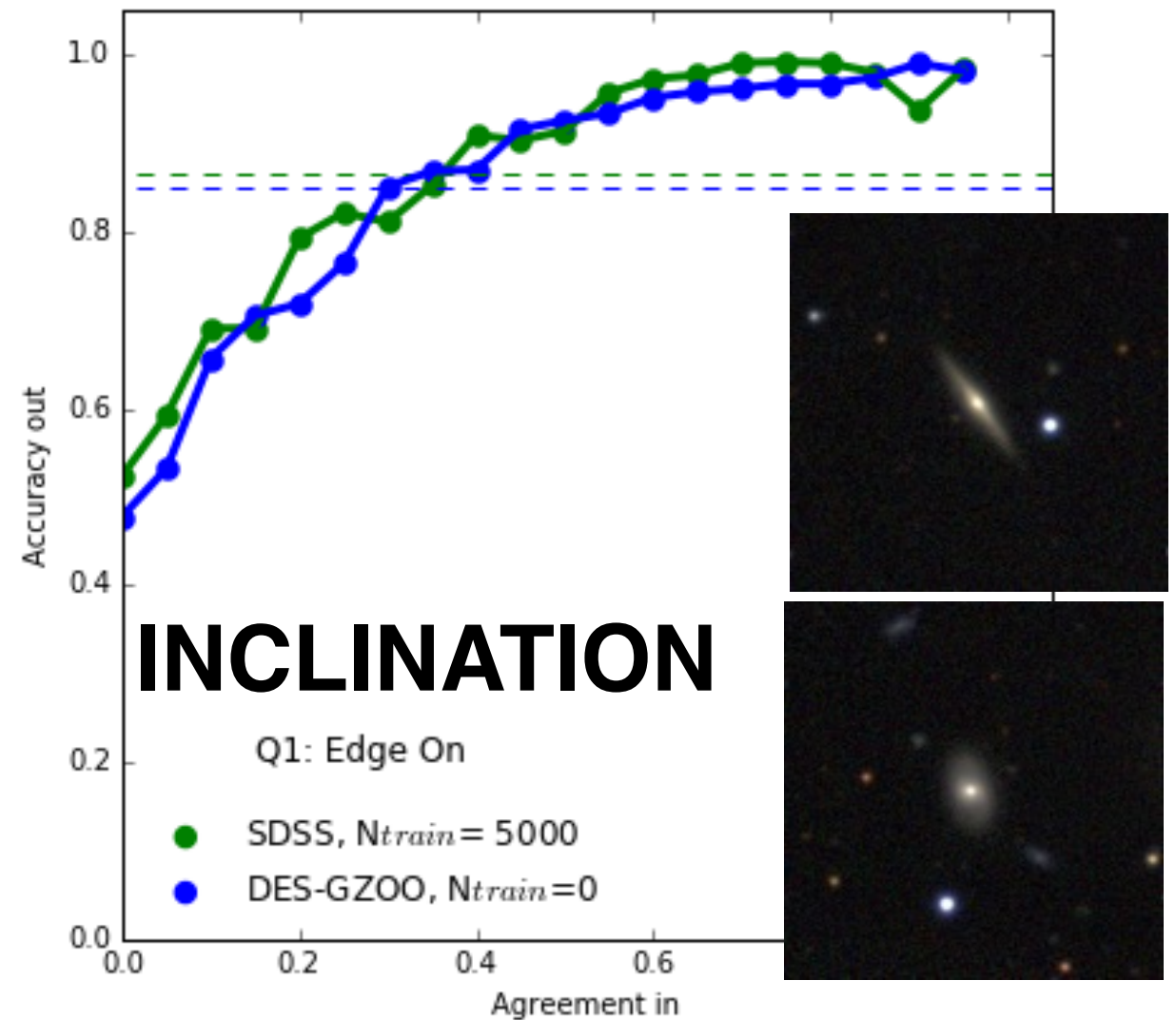
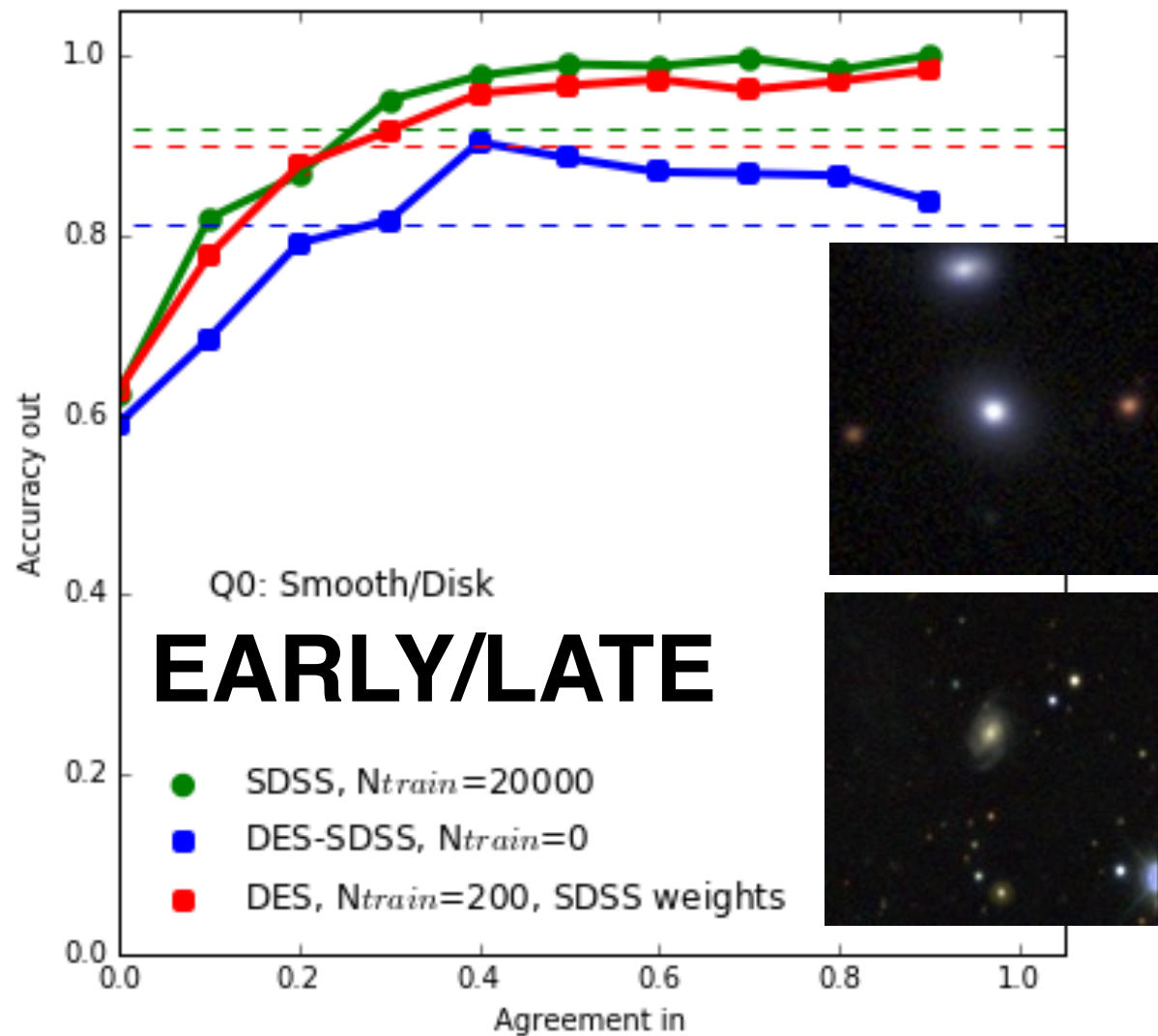


**DES**



# Knowledge transfer from SDSS to DES

DOMINGUEZ-SANCHEZ, HUERTAS-COMPANY, BERNARDI et al. 17b

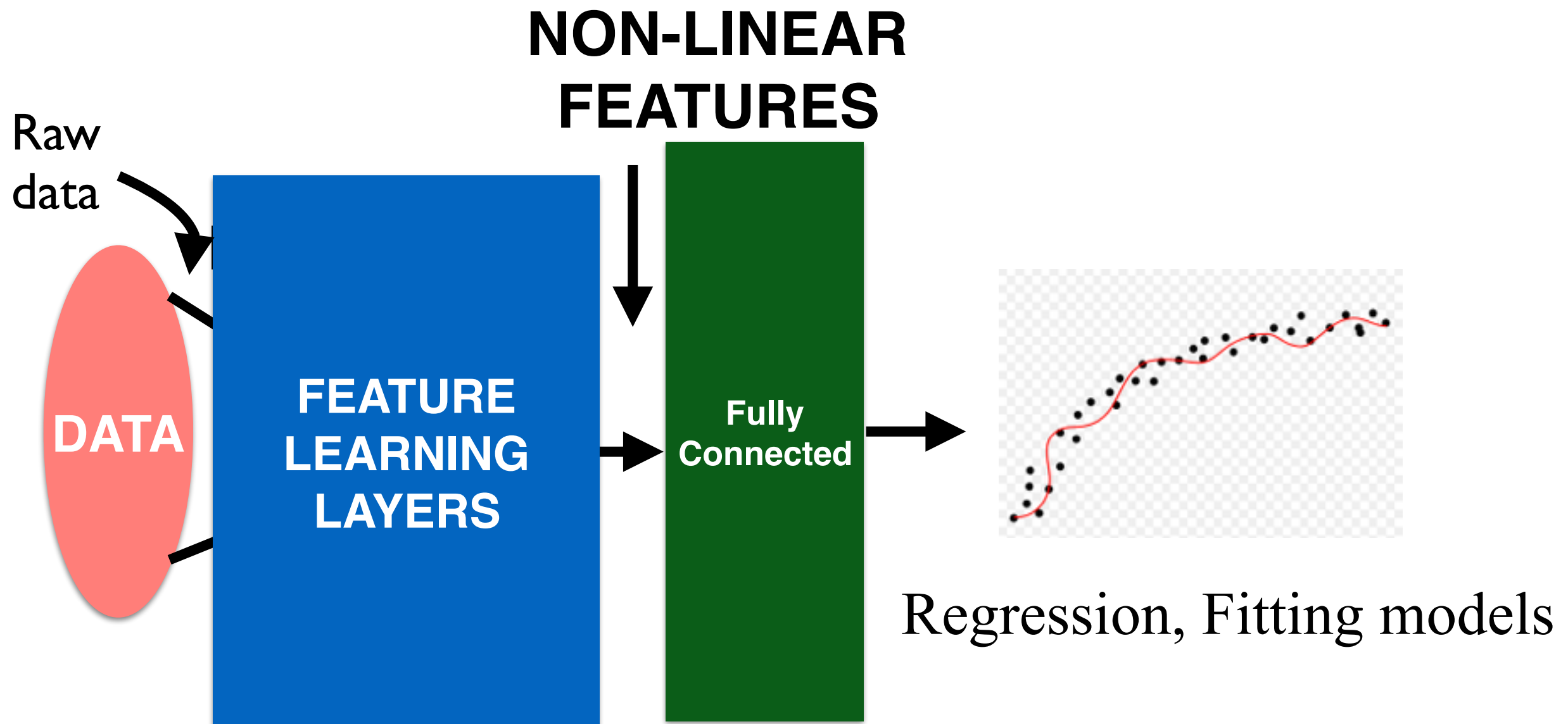


Only 200 (1%!) objects classified in DES are needed to reach an accuracy  $>90\%$  if a machine trained on the SDSS is used

For some properties, i.e. EDGE-ON galaxies. No training at all is needed to go from SDSS to DES



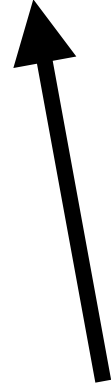
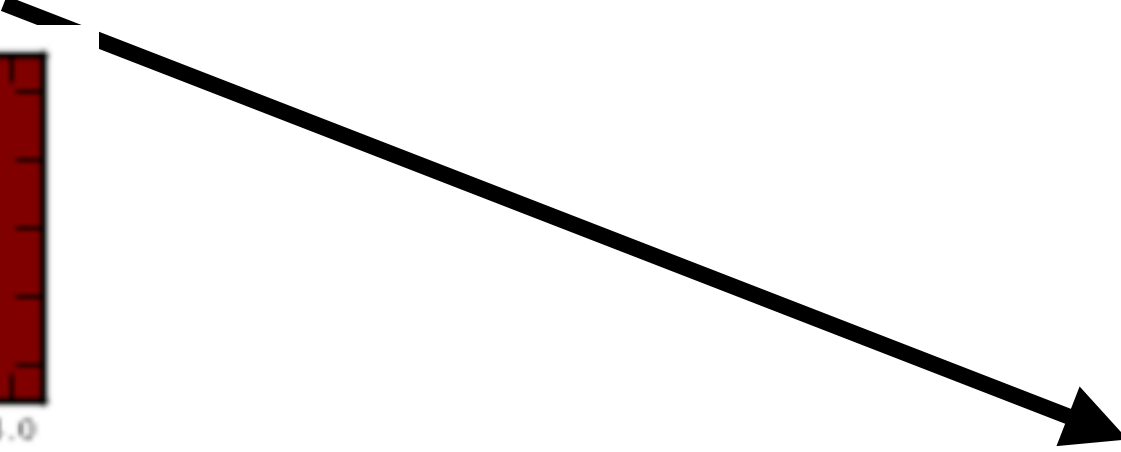
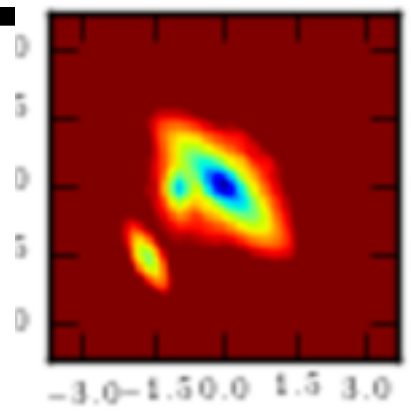
**GROUP #2:** Efficient and fast quantitative measurements on large amount of (multi-lambda) data [photoz's, sizes, ellipticities]



# MORPHOMETRICS



flux, size, axis-ratio, PA,  
Sersic index  
*[PSF corrected]*



**MODEL**



**MODEL\*PSF**

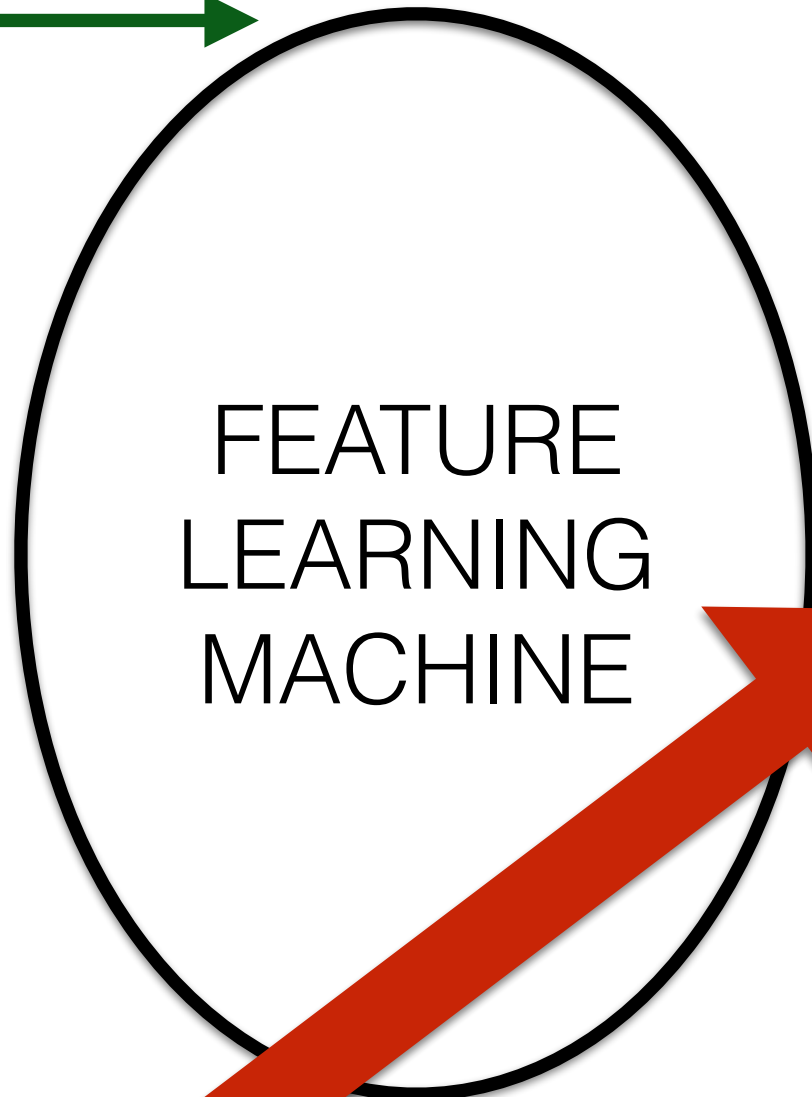


**MINIMIZATION**

**TRAINING:**

simulations of  
analytic profiles  
with PSF, noise  
effects

(no limits on the size)



Flux



Sersic Index



Radii



b/a

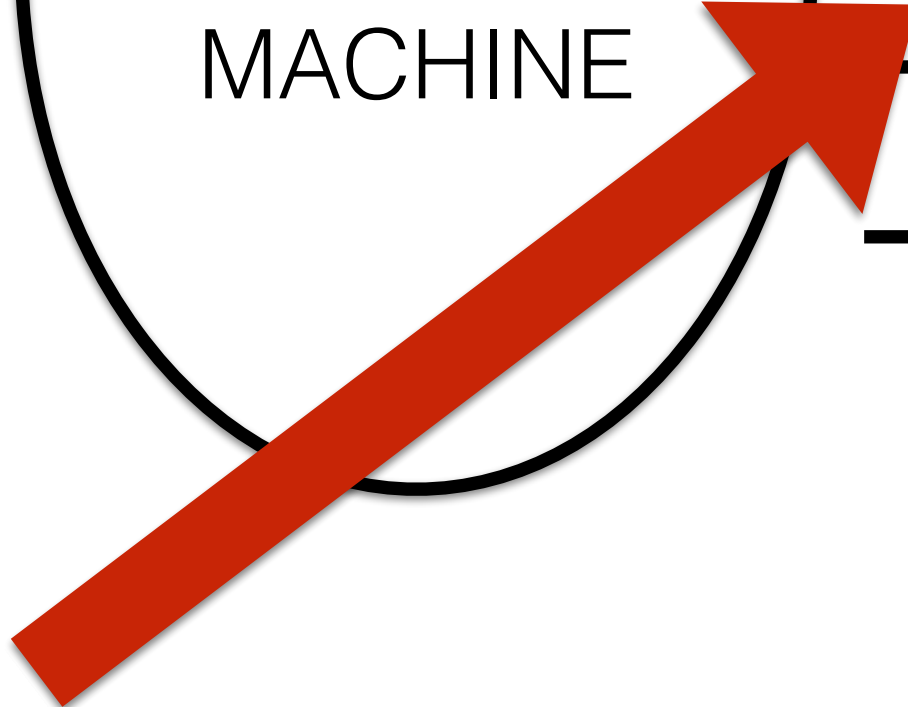
.

.

.

**DATA:**

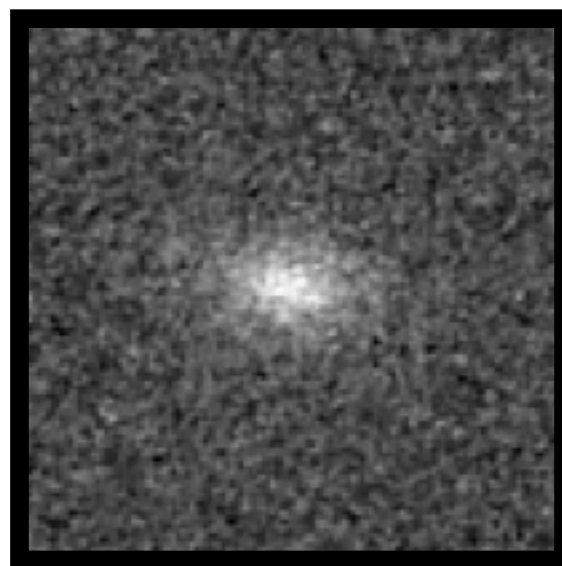
HST deep field  
observations  
CANDELS

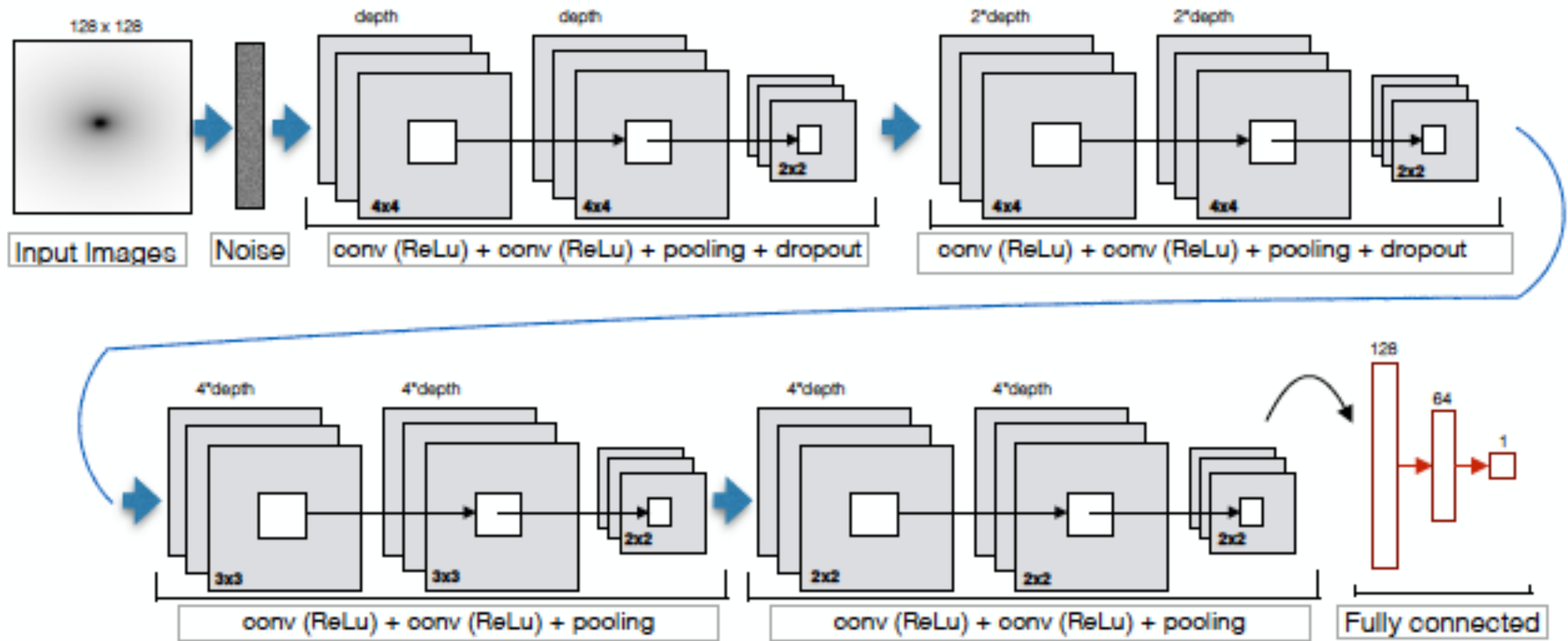




# Standard analytic profiles

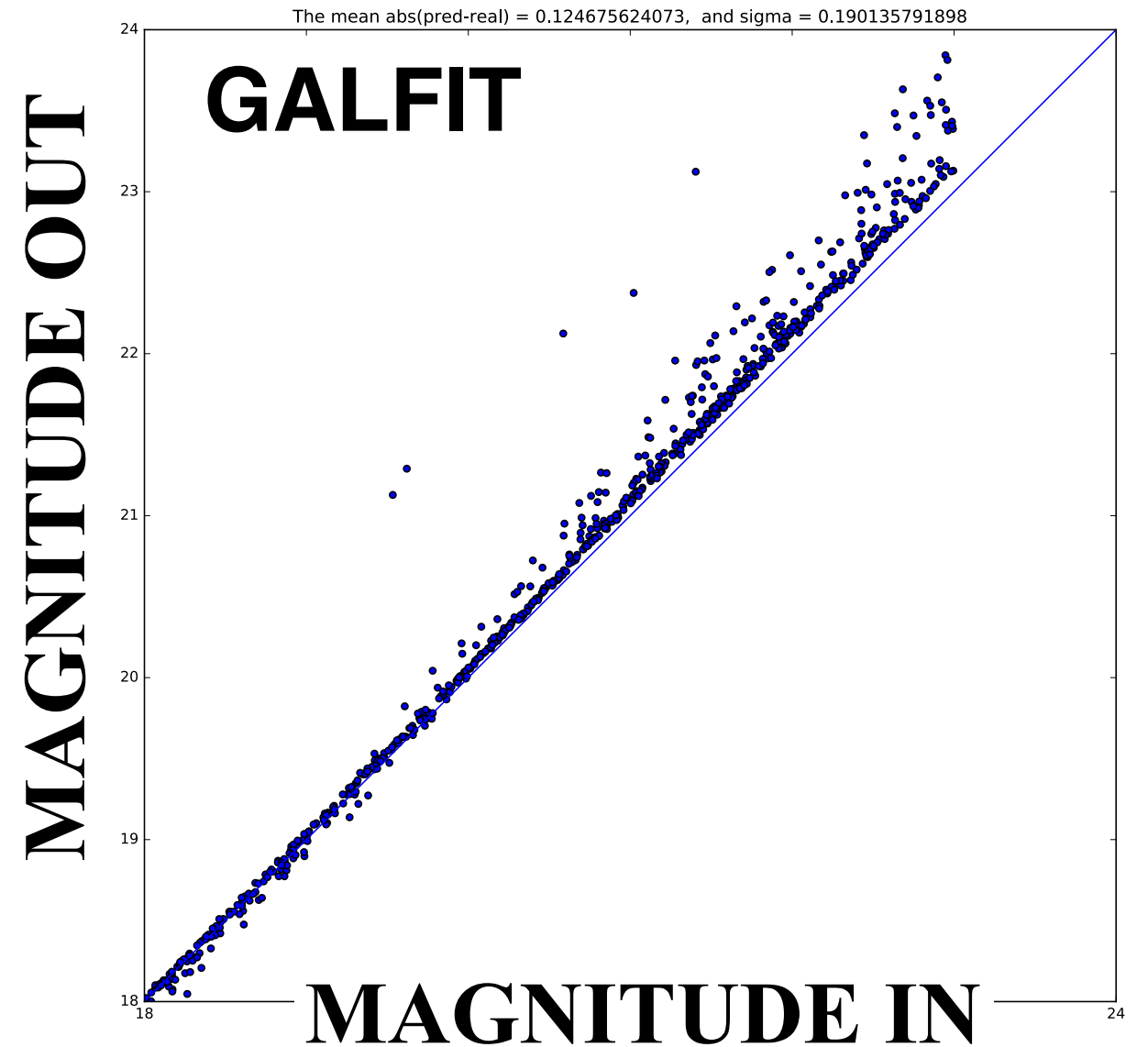
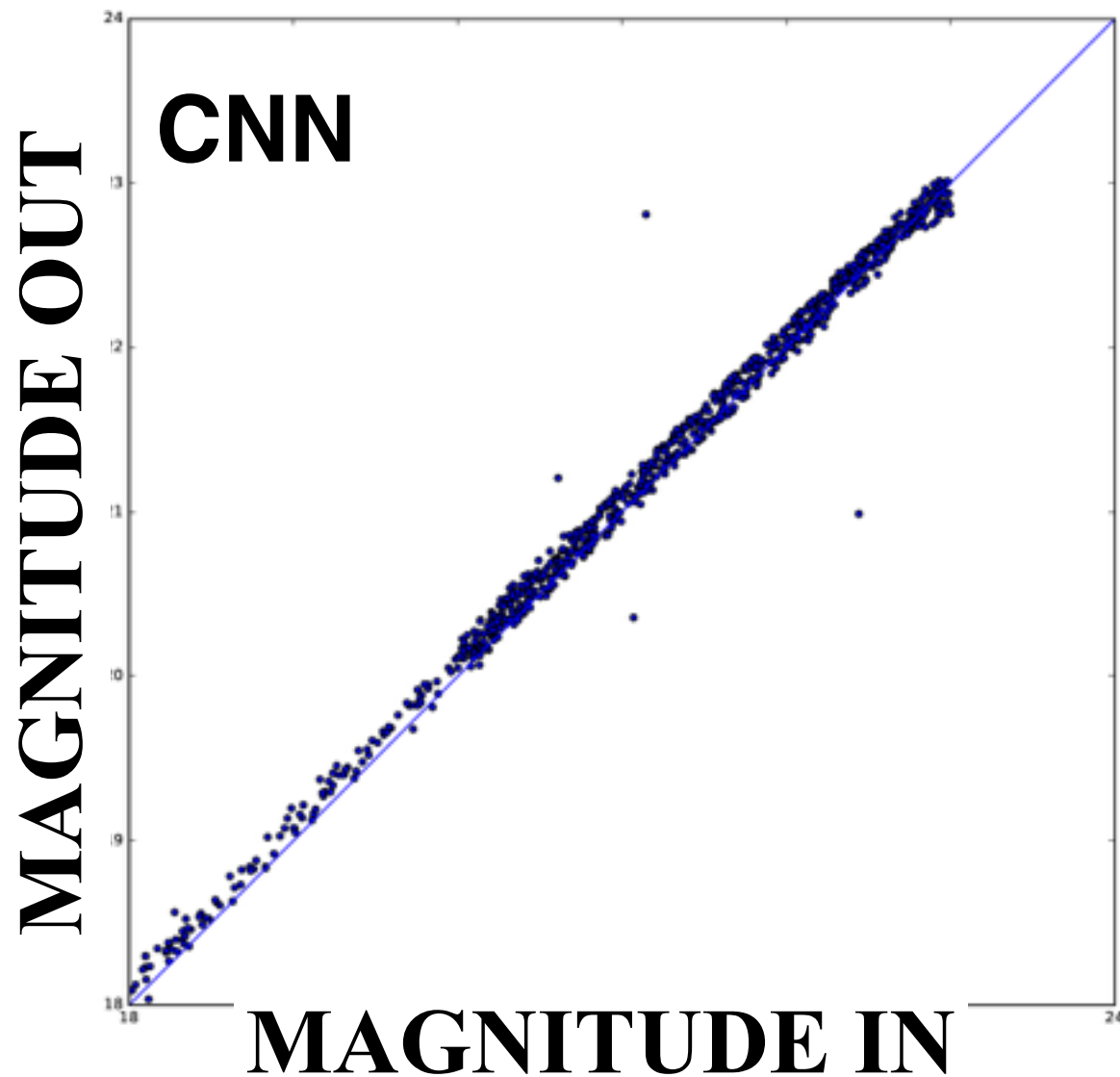
- 100.000-300.000 galaxies [GALSIM]
  - Real HST background added + PSF (F160)
  - Random distribution of parameters (uniform):
    - $18 < \text{Mag} < 24$ ,  $0 < \text{BT} < 1$ ,  $< \text{Nb} <$ ,  $\text{Nd} = 1$ ,  $0.2 < \log(\text{rb}) < 1.3$ ,  
 $0.2 < \log(\text{rd}) < 1.5$ ,  $0.05 < \text{eb} < 0.95$ ,  $0.05 < \text{ed} < 0.95$ ,  $0 < \text{PA} < 180$
  - $64 * 64$  stamps
  - **FULLY IDEALISTIC -**  
**NO COMPANIONS**  
**NO IRREGULARS**  
**NO CLUMPY!**





# ON SIMULATIONS

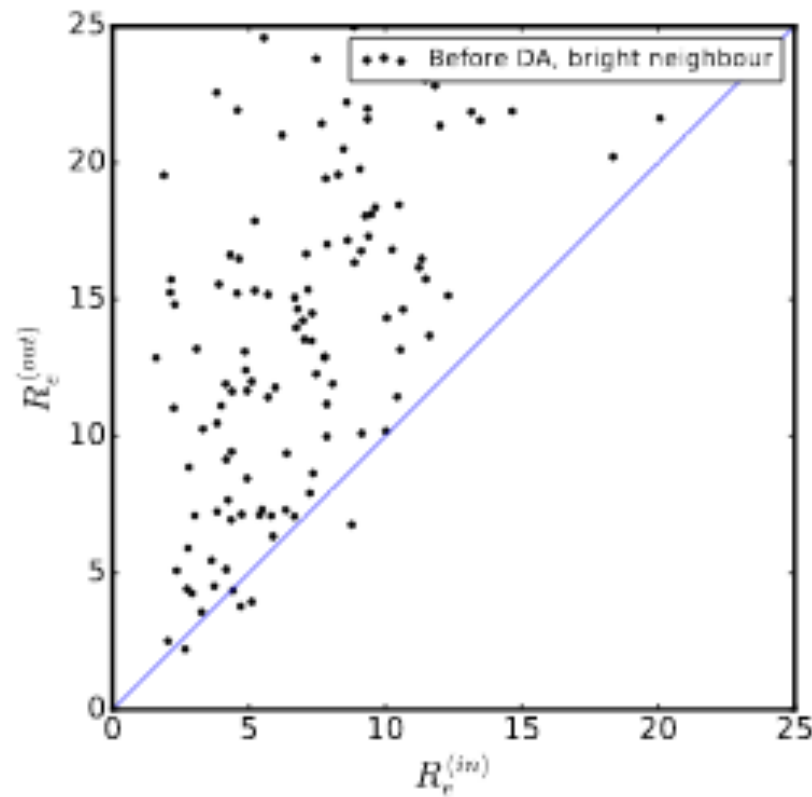
TUCCILLO, HUERTAS-COMPANY et al. 17



VERY SIMILAR RESULTS ON THE SAME SIMULATIONS, BUT  
CNNs are several orders of magnitude faster [**3.5 hrs vs. <1 sec for  
~1000 objects**]

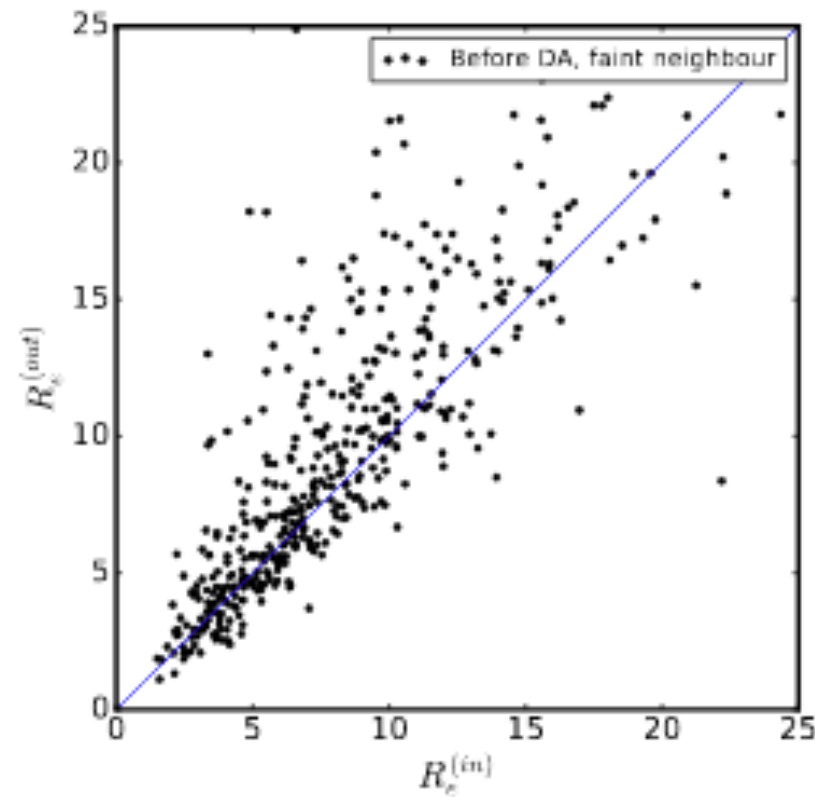
# MORPHOMETRY OF REAL GALAXIES TRAINED ON ANALYTIC PROFILES

## BRIGHT NEIGH.



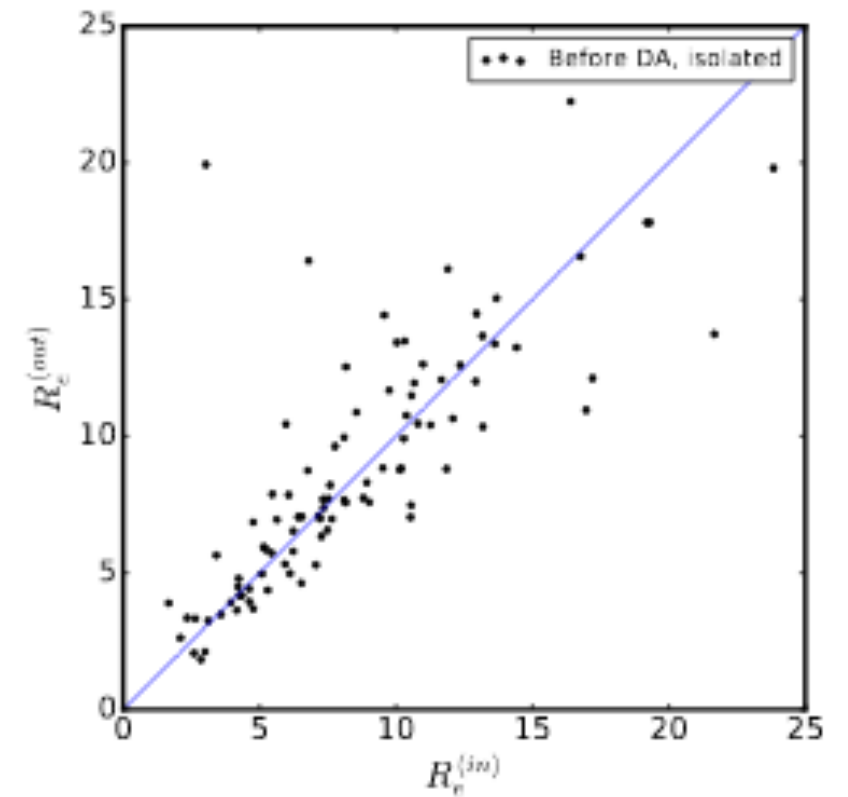
(c) BDA bright neighbours

## FAINT NEIGH.



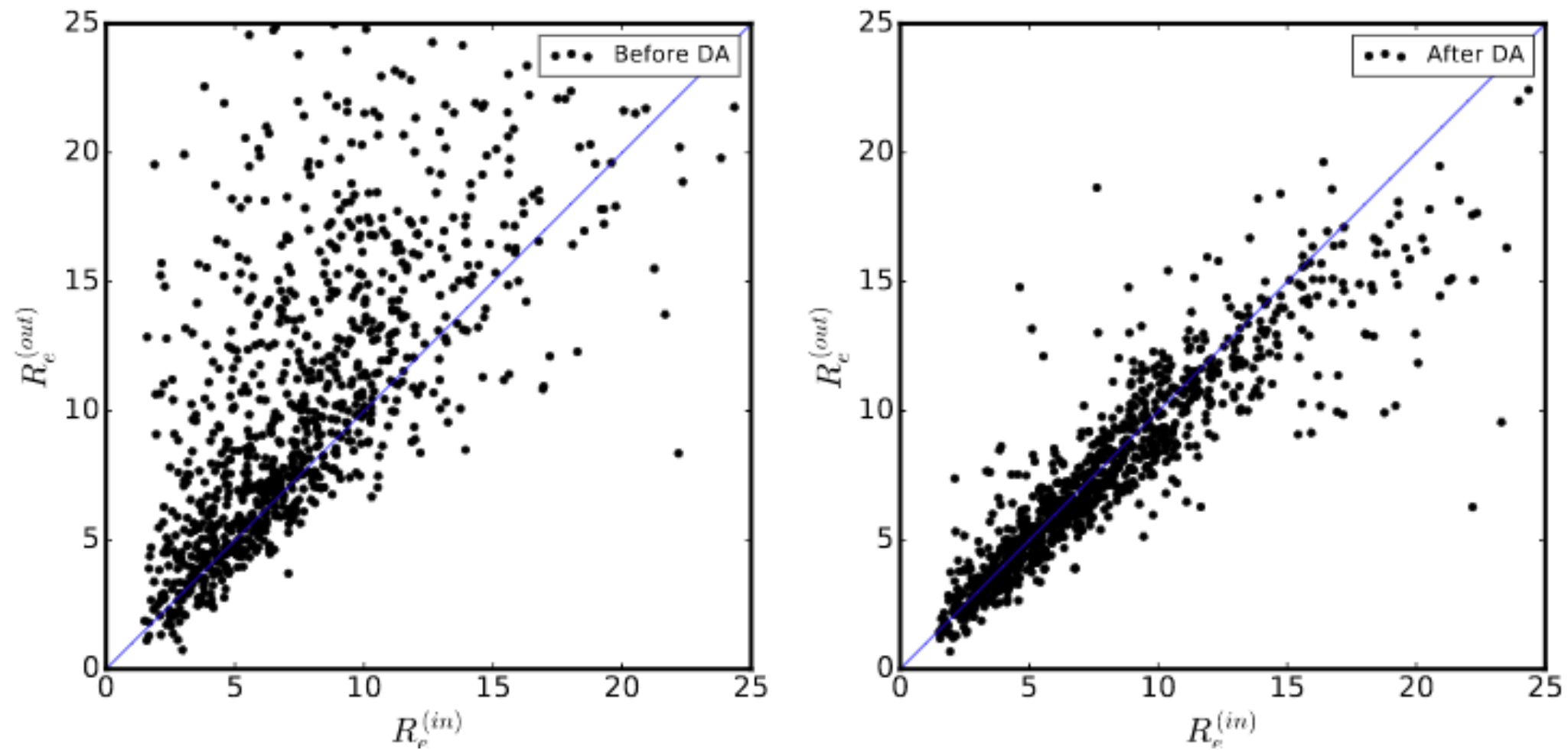
(d) BDA faint neighbours

## ISOLATED

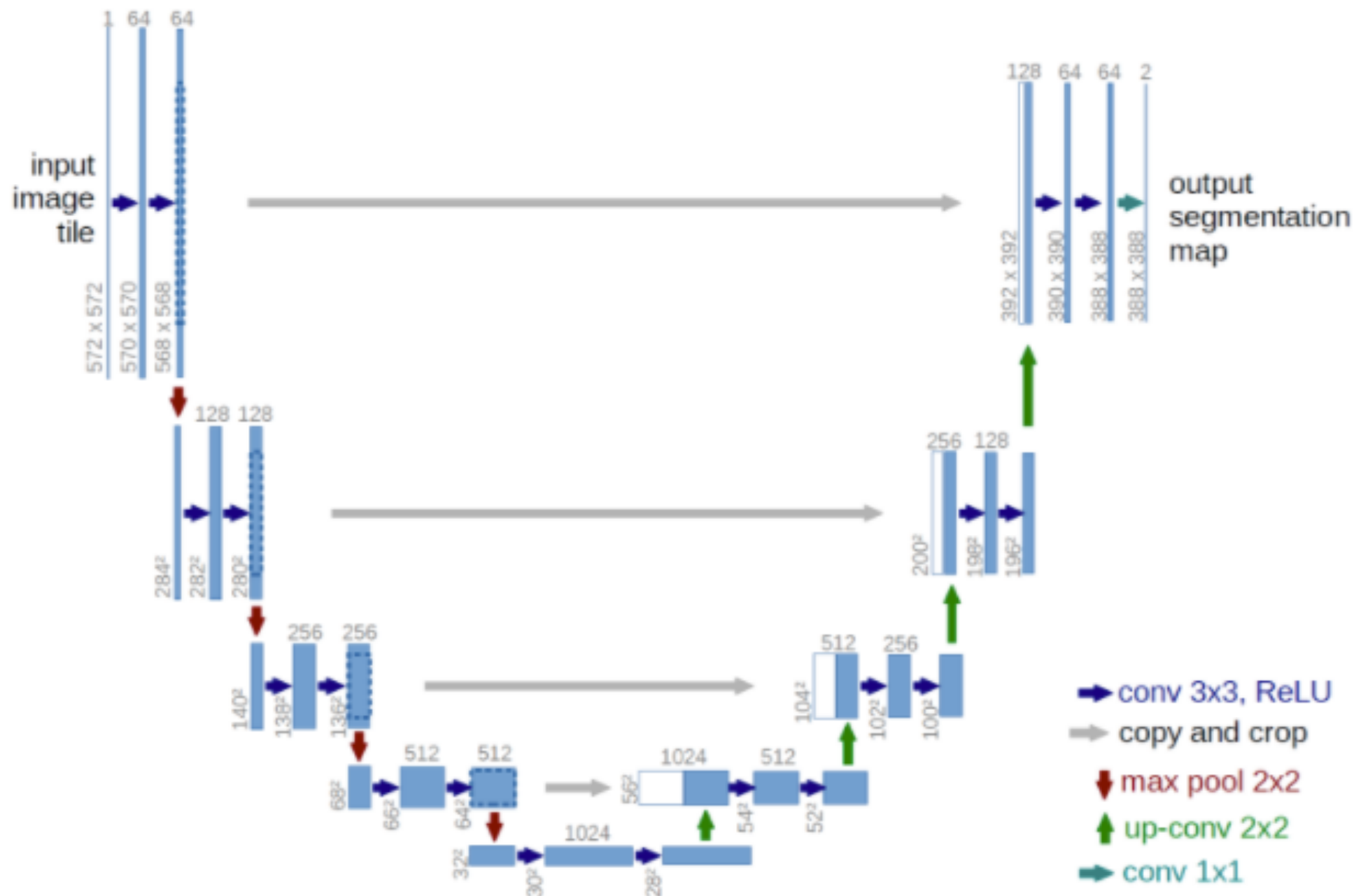


(e) BDA isolated galaxies

# DOMAIN ADAPTATION: 0.1% OF “REAL” GALAXIES



# Coming soon: U-net for bulge/disc decompositions...



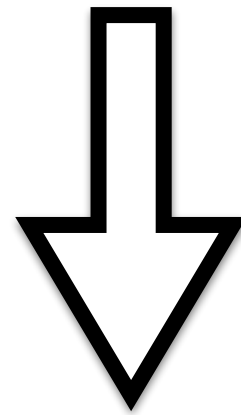


# VELA hydrodynamic simulations

**Ceverino+15**

35 high res ( $\sim 20\text{pc}$ ) zoom-in simulations  
hydroART

radiative and supernovae feedback  
stops at  $z=1$  -  $M_h=10^{11}-2\cdot 10^{12}$

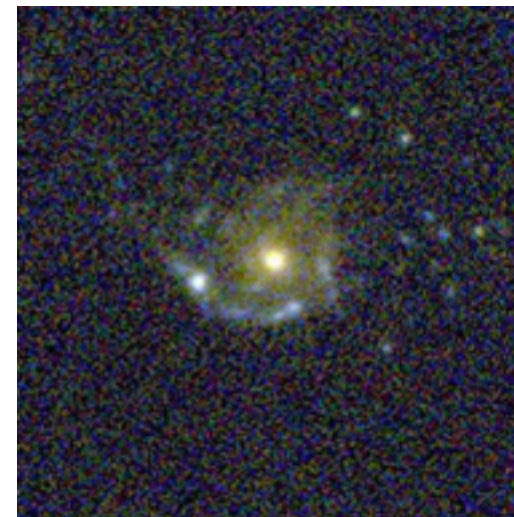
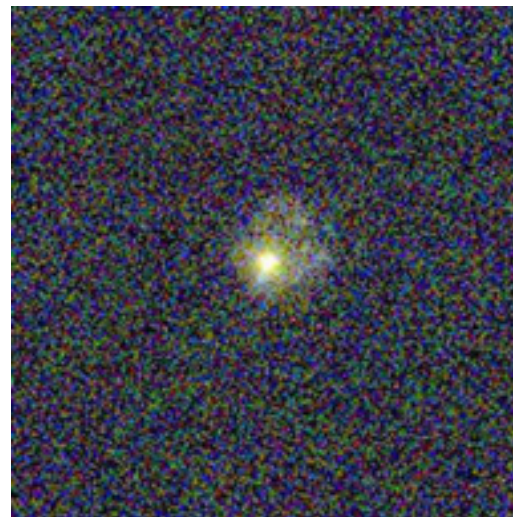
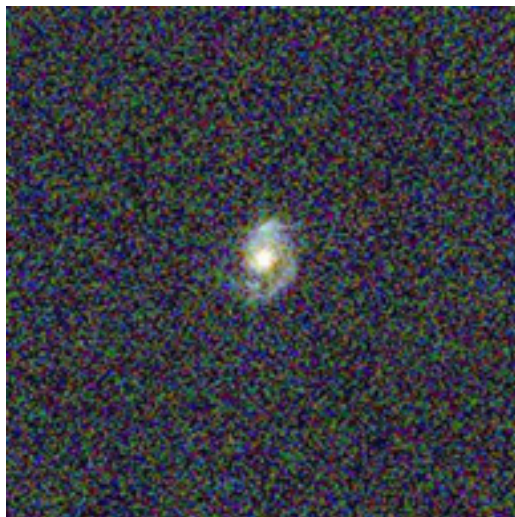


mock images [sunrise]

$T_{\text{step}} \sim 200\text{Myrs}$

10 projections

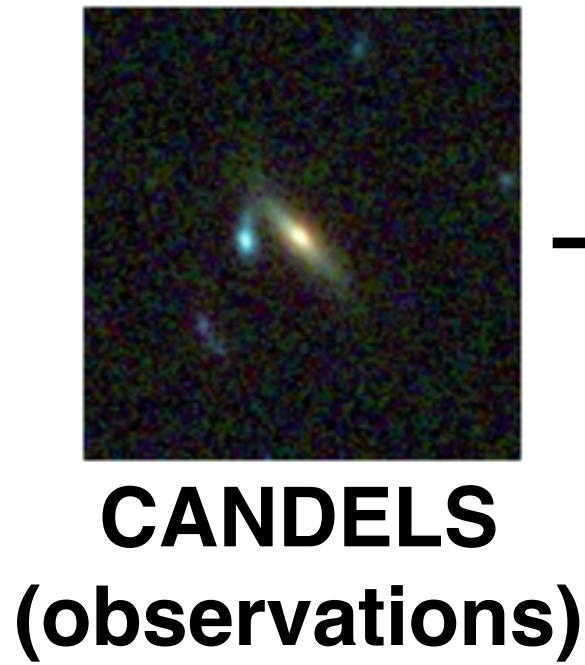
HST like



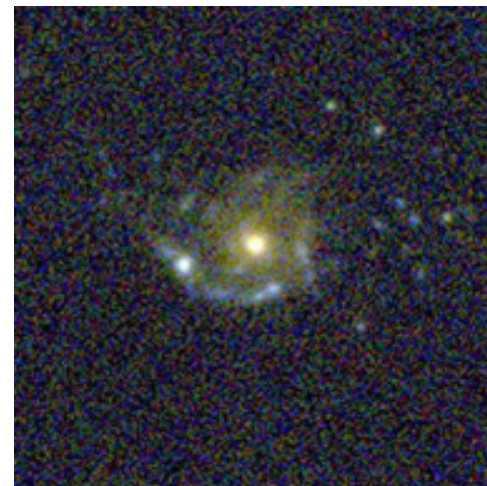
**[G. Snyder, J. Lotz]**

**Margalef, MHC in prep.**

DEEPLLEGATO



Re (2D, light)



Re (2D, light)

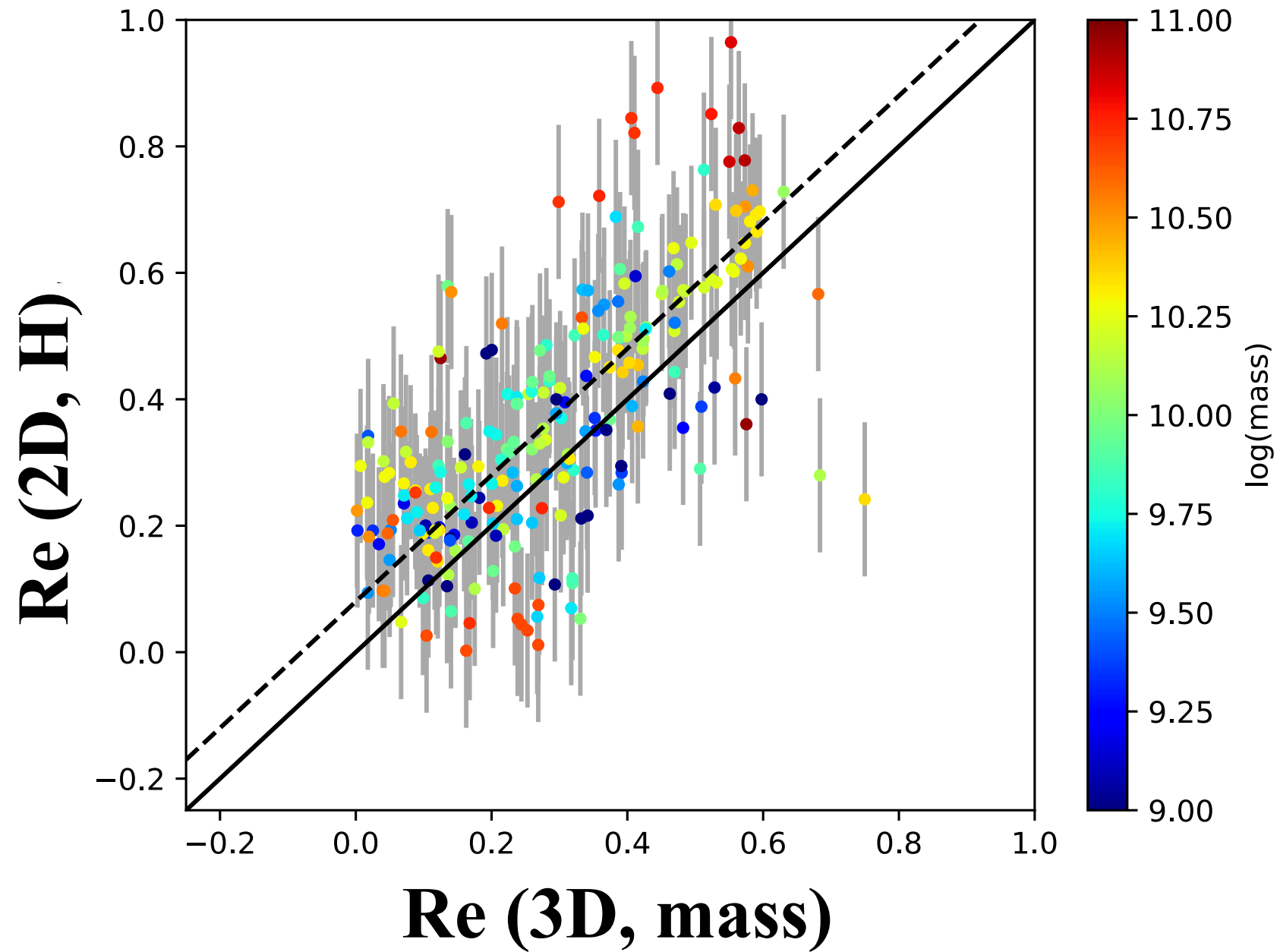
**VELA  
(simulation)**

SIMULATION  
METADATA

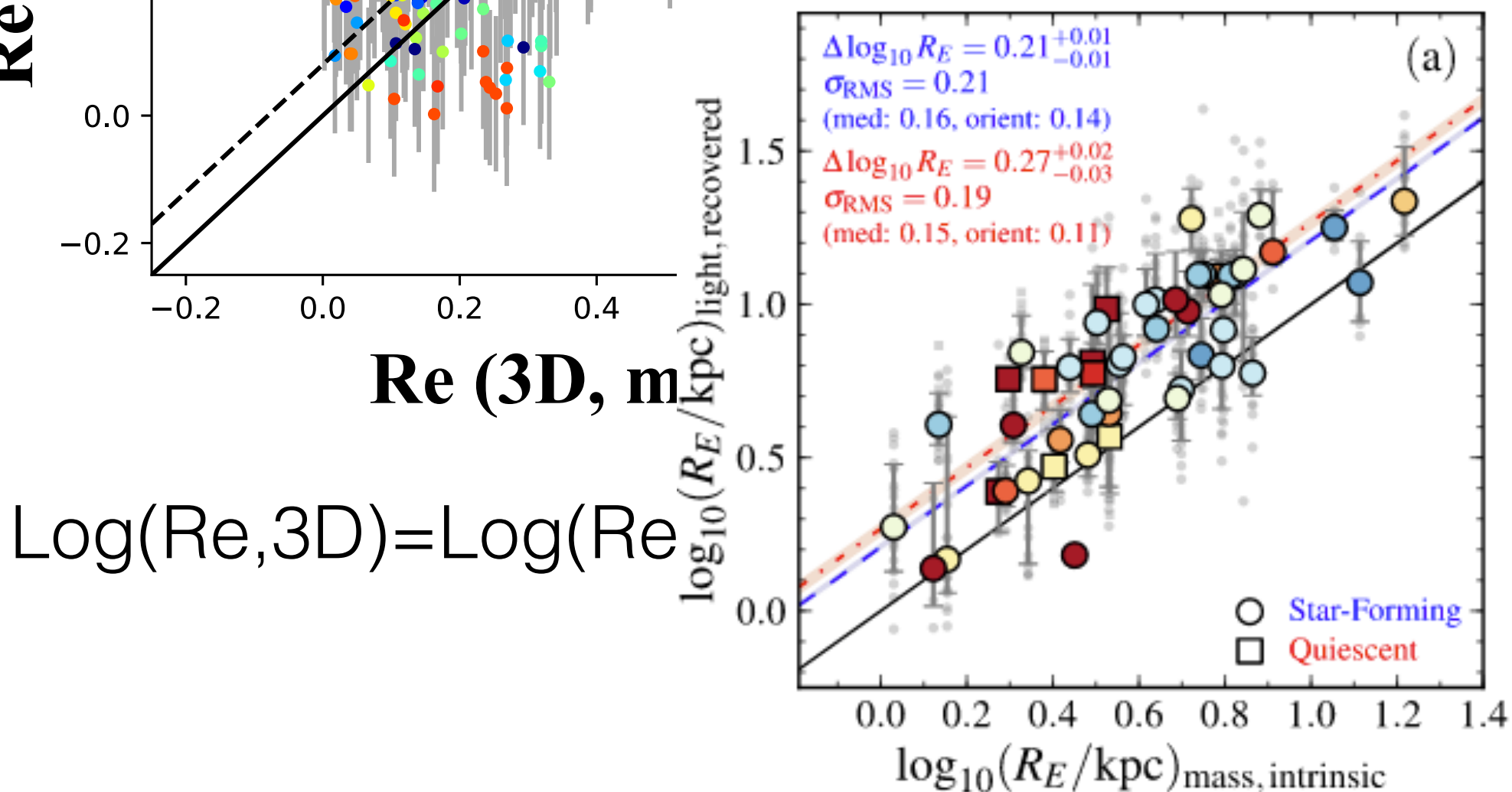
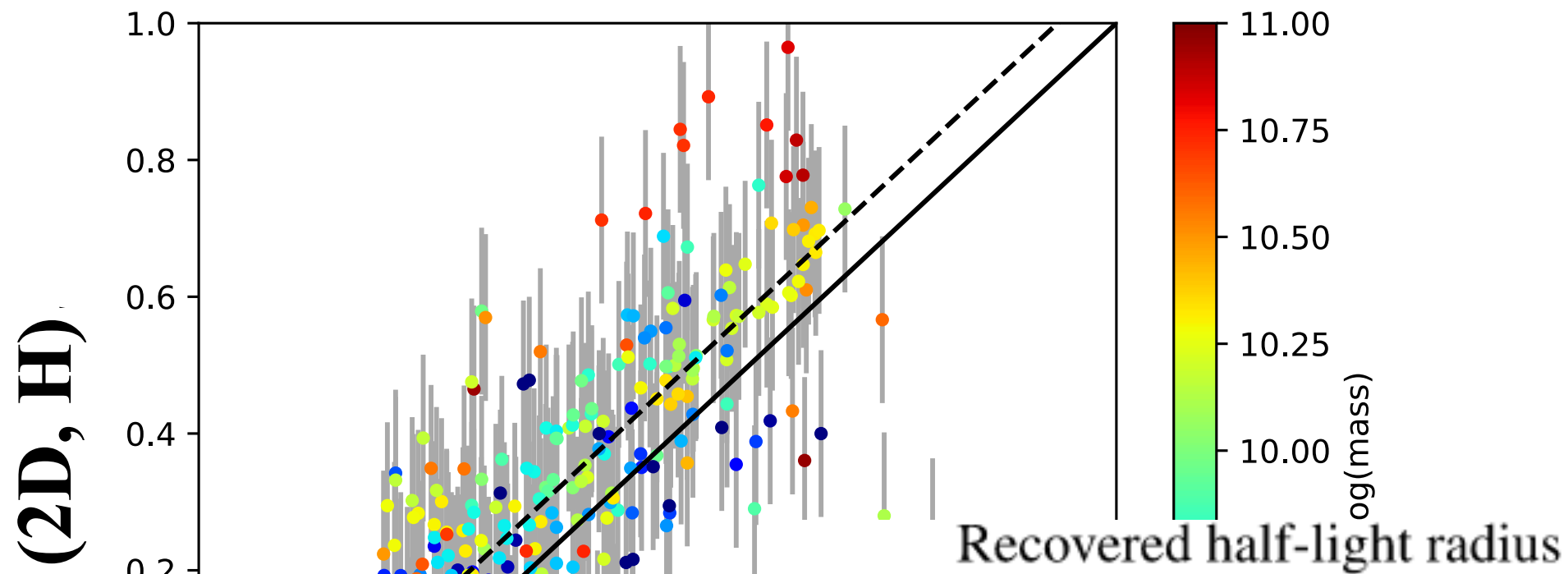
Re (3D, mass,  $<0.1R_{\text{vir}}$ )



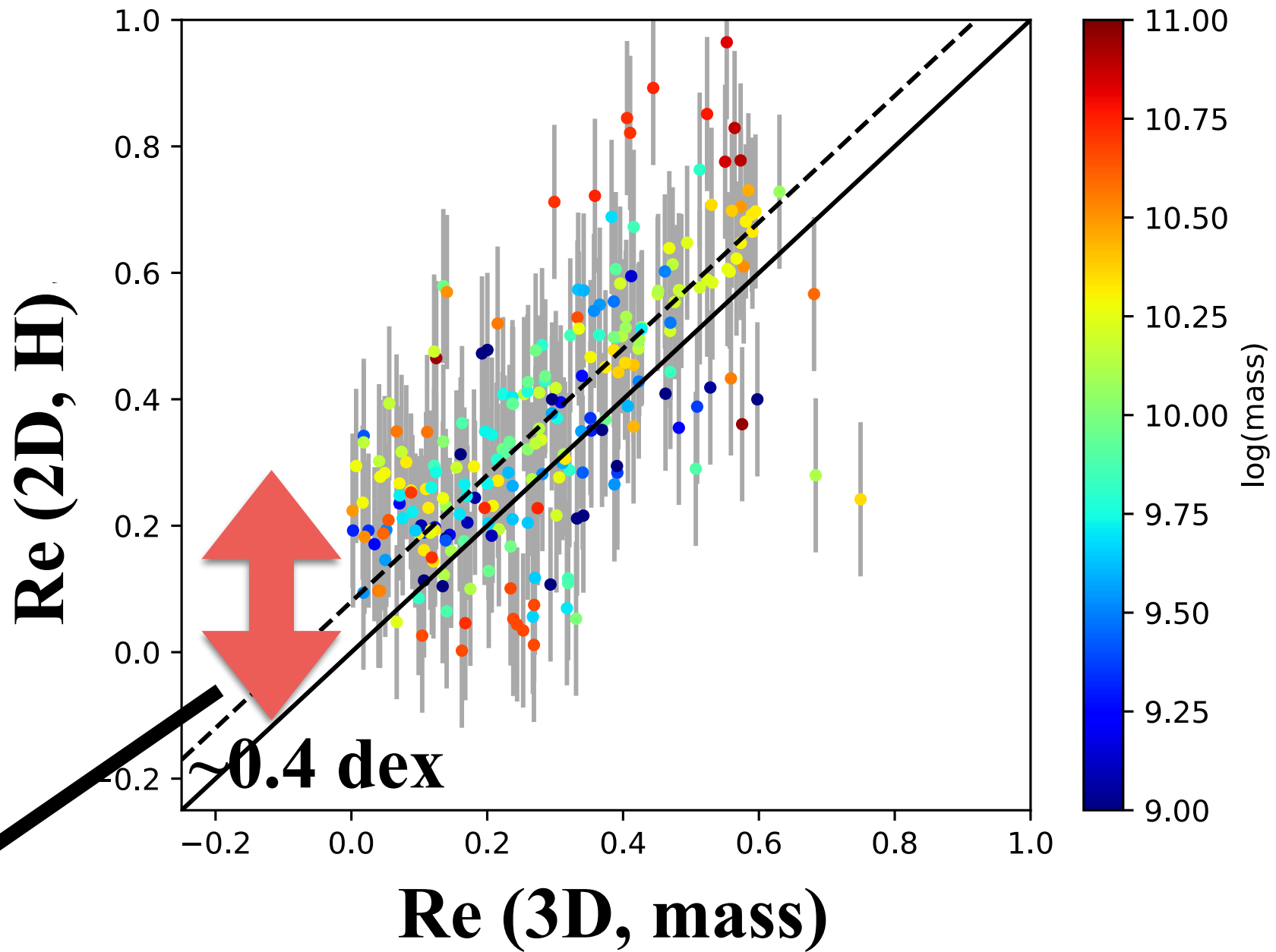




$$\text{Log}(\text{Re}, 3\text{D}) = \text{Log}(\text{Re}, 2\text{D}) * 0.8 + 0.13$$



$\text{Log}(Re, 3D) = \text{Log}(Re$



**reduce  
scatter  
with DL?**

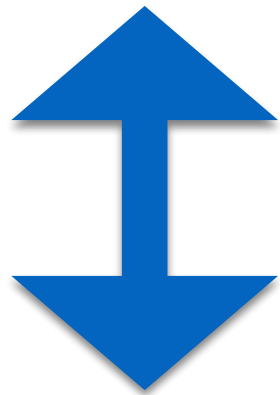
$$\text{Log(Re,3D)} = \text{Log(Re,2D)} * 0.8 + 0.13$$

- **GROUP #3:** Find new hidden observables in the data, - Linking observations and theory

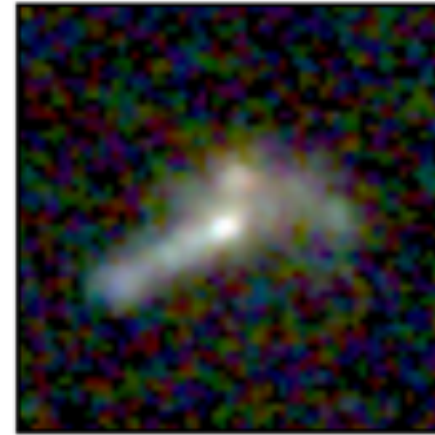
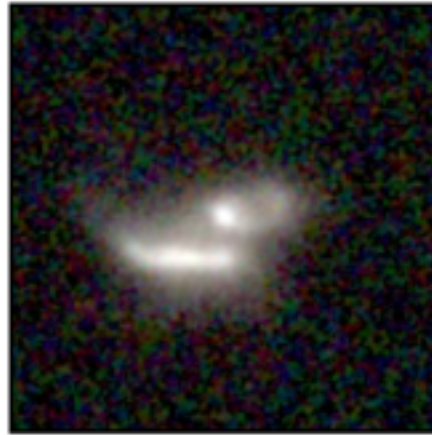
# HYDRODYNAMIC SIMULATIONS

(e.g. Horizon-AGN, Illustris,  
FIRE, VELA...)

ASSEMBLY  
PROCESSES



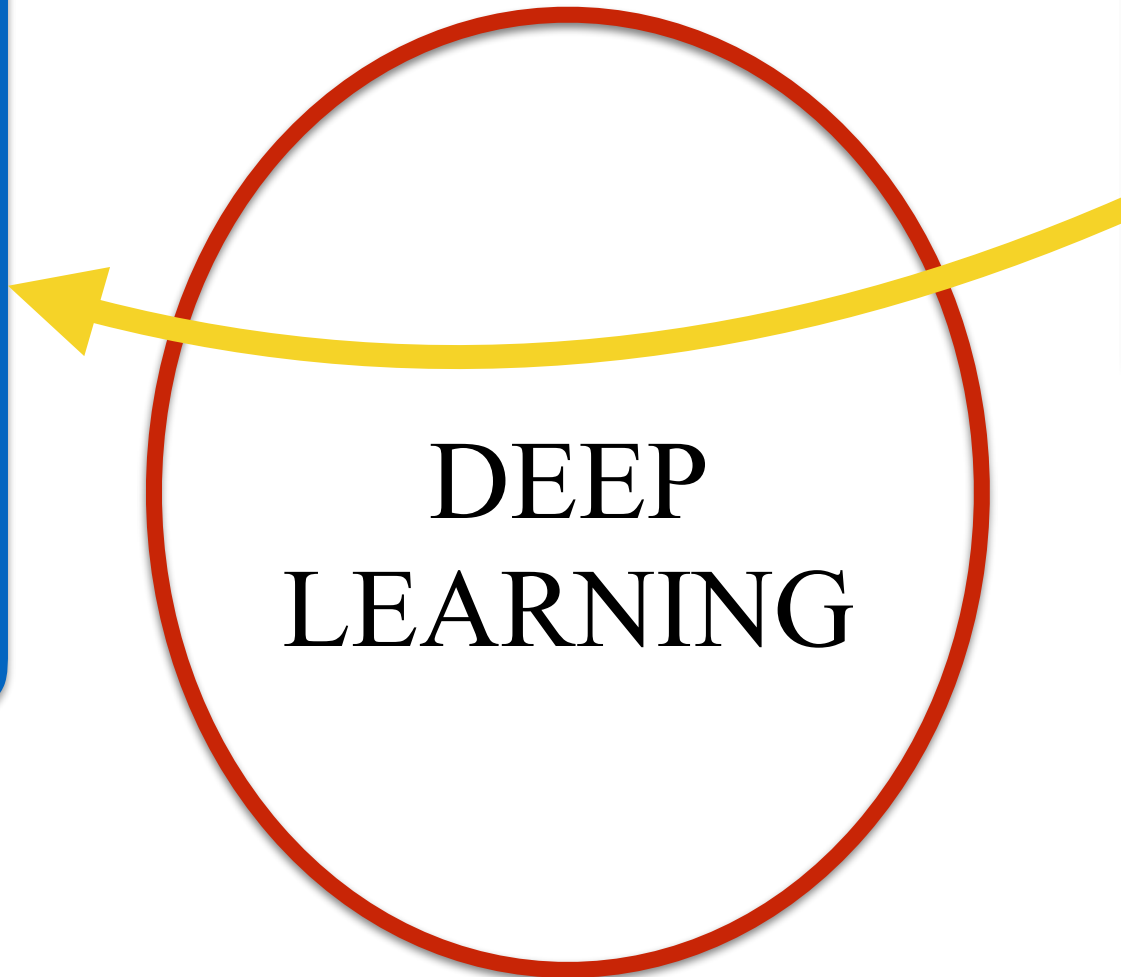
MOCK IMAGES



## DATA

[OBSERVATIONS]

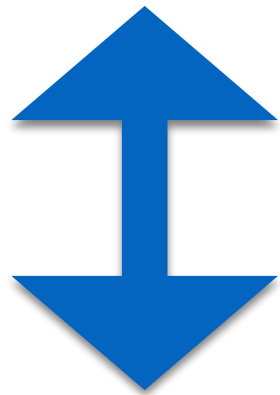
DEEP  
LEARNING



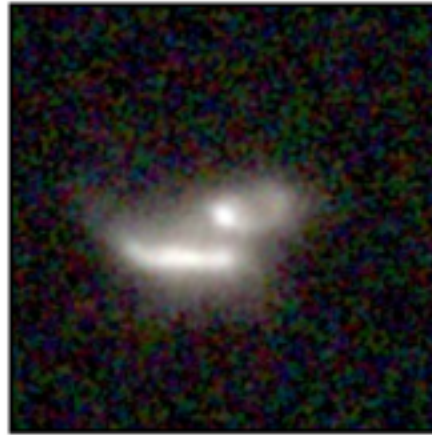
# HYDRODYNAMIC SIMULATIONS

(e.g. Horizon-AGN, Illustris,  
FIRE, VELA...)

ASSEMBLY  
PROCESSES



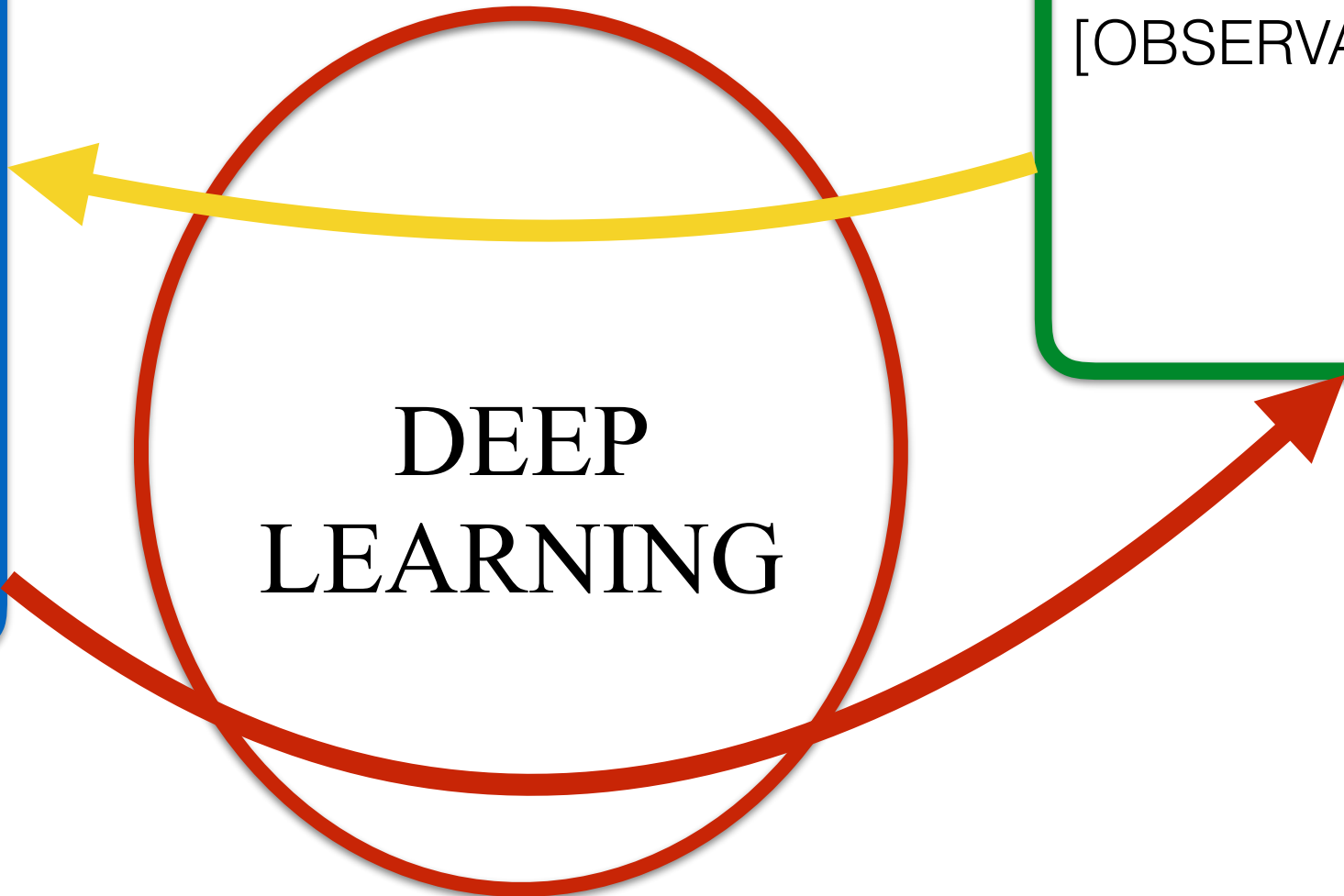
MOCK IMAGES



## DATA

[OBSERVATIONS]

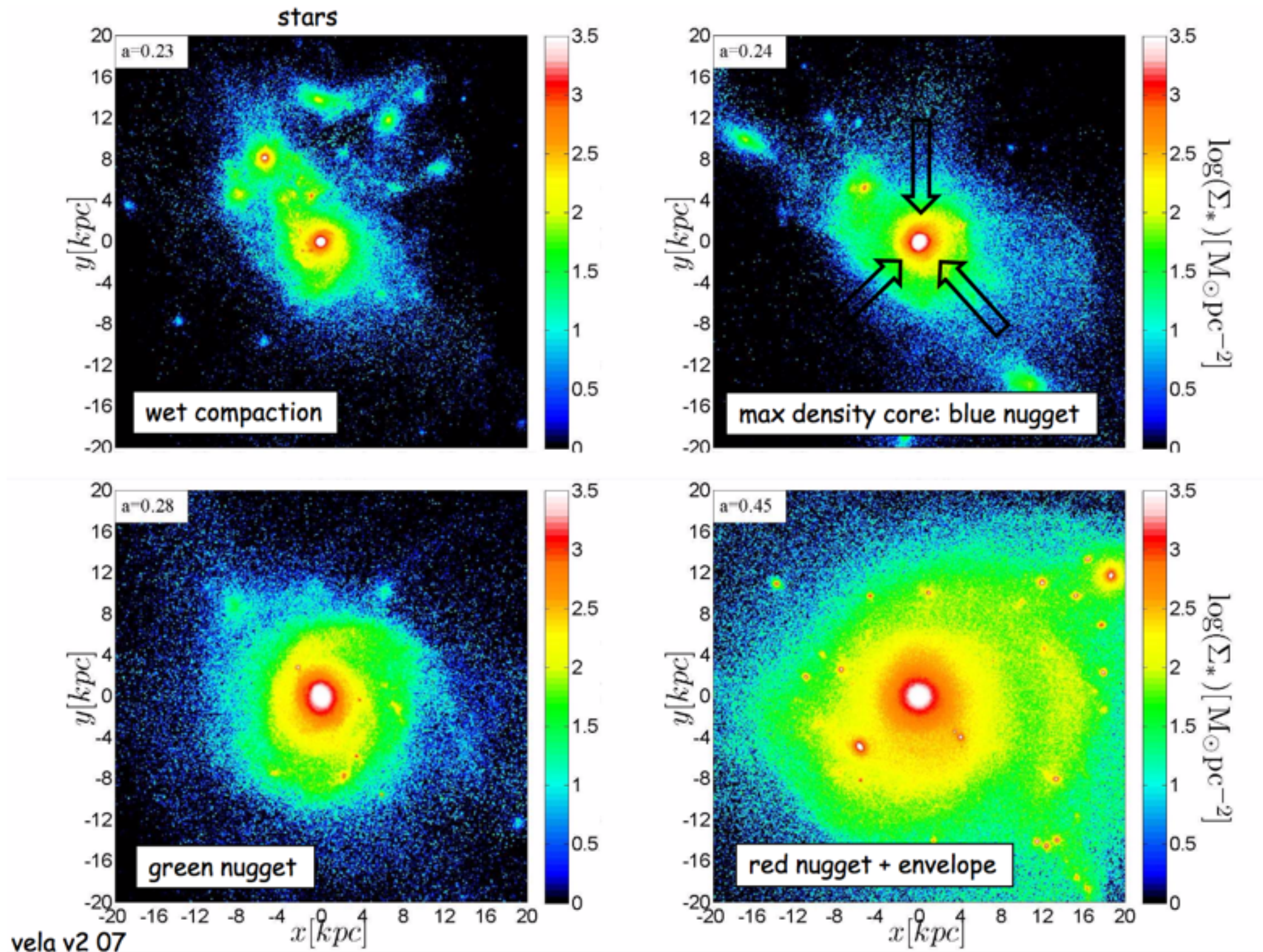
DEEP  
LEARNING



Is “wet compaction” a common mechanism for bulge formation?



# Is “wet compaction” a common mechanism for bulge formation?

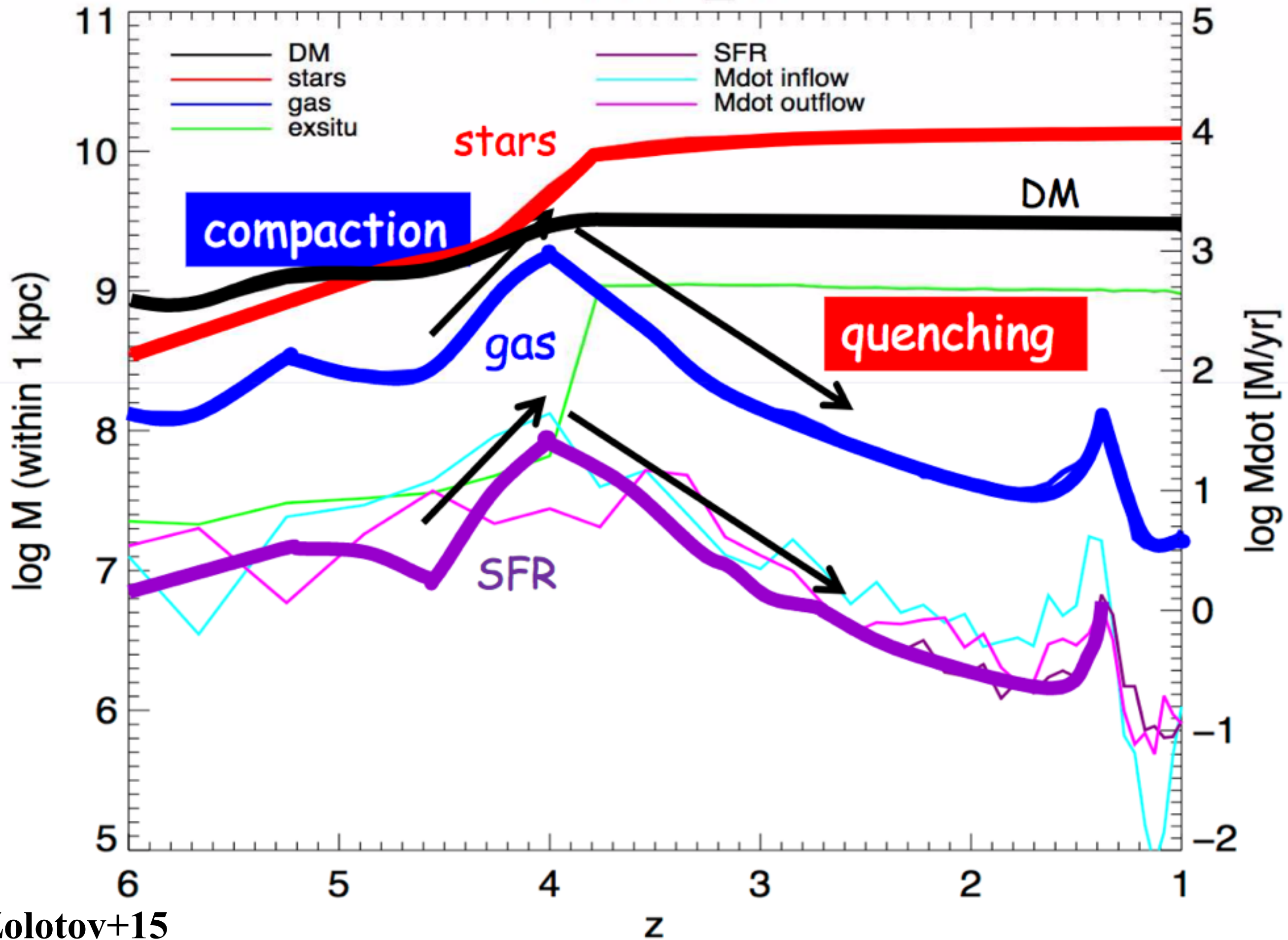


**Ceverino+15**  
**Zolotov+15**  
**Tacchella+17**

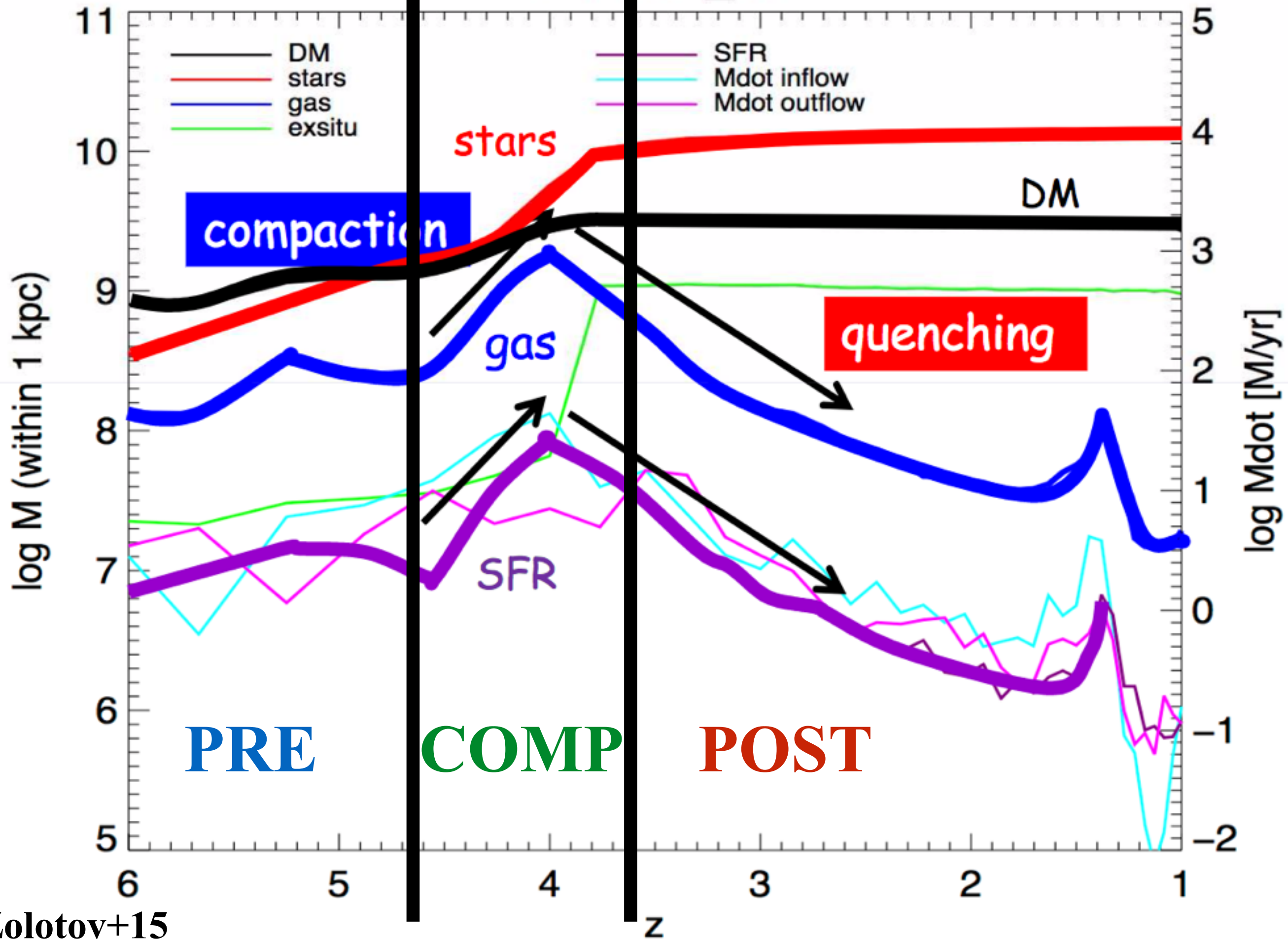
Courtesy of A. Dekel

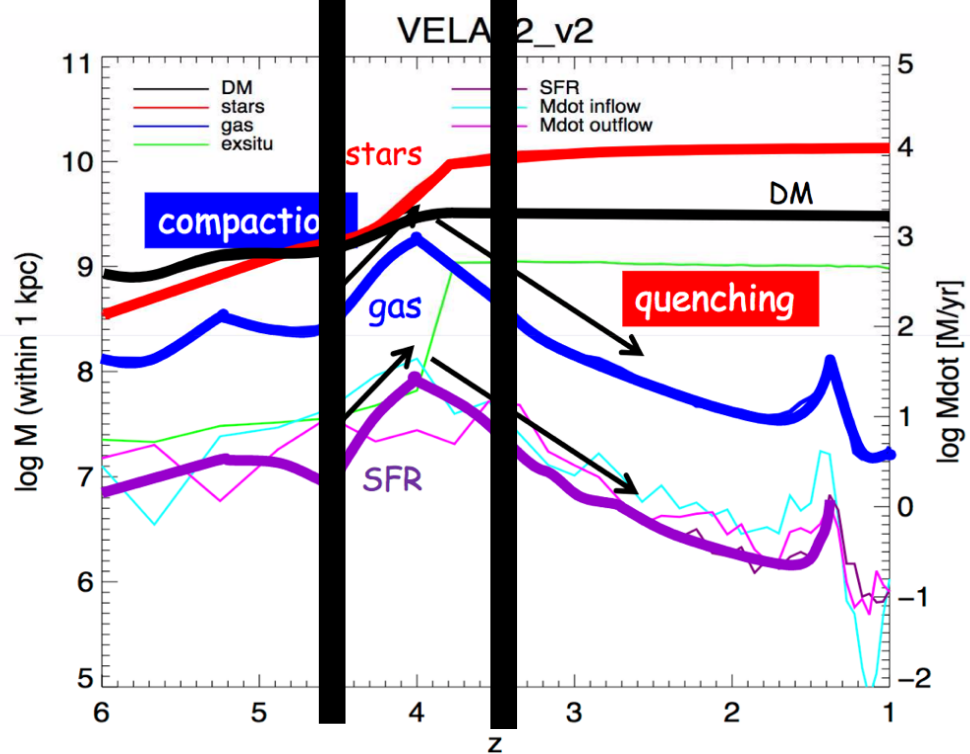


# VELA12\_v2



# VELA12\_v2

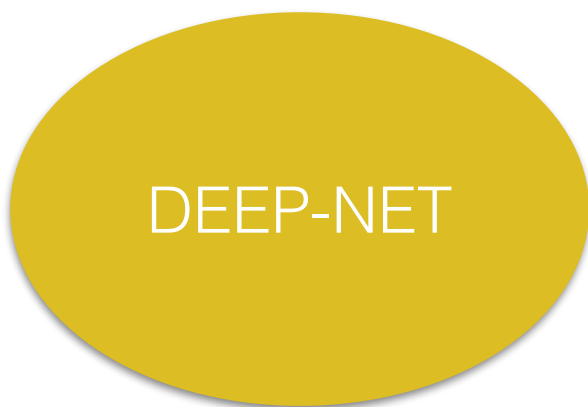




**“WET COMPACTION”**  
**[simulation metadata]**

ANY SIGNATURE  
 OF  
 “COMPACTION”  
 IN THE STELLAR  
 DISTRIBUTION?

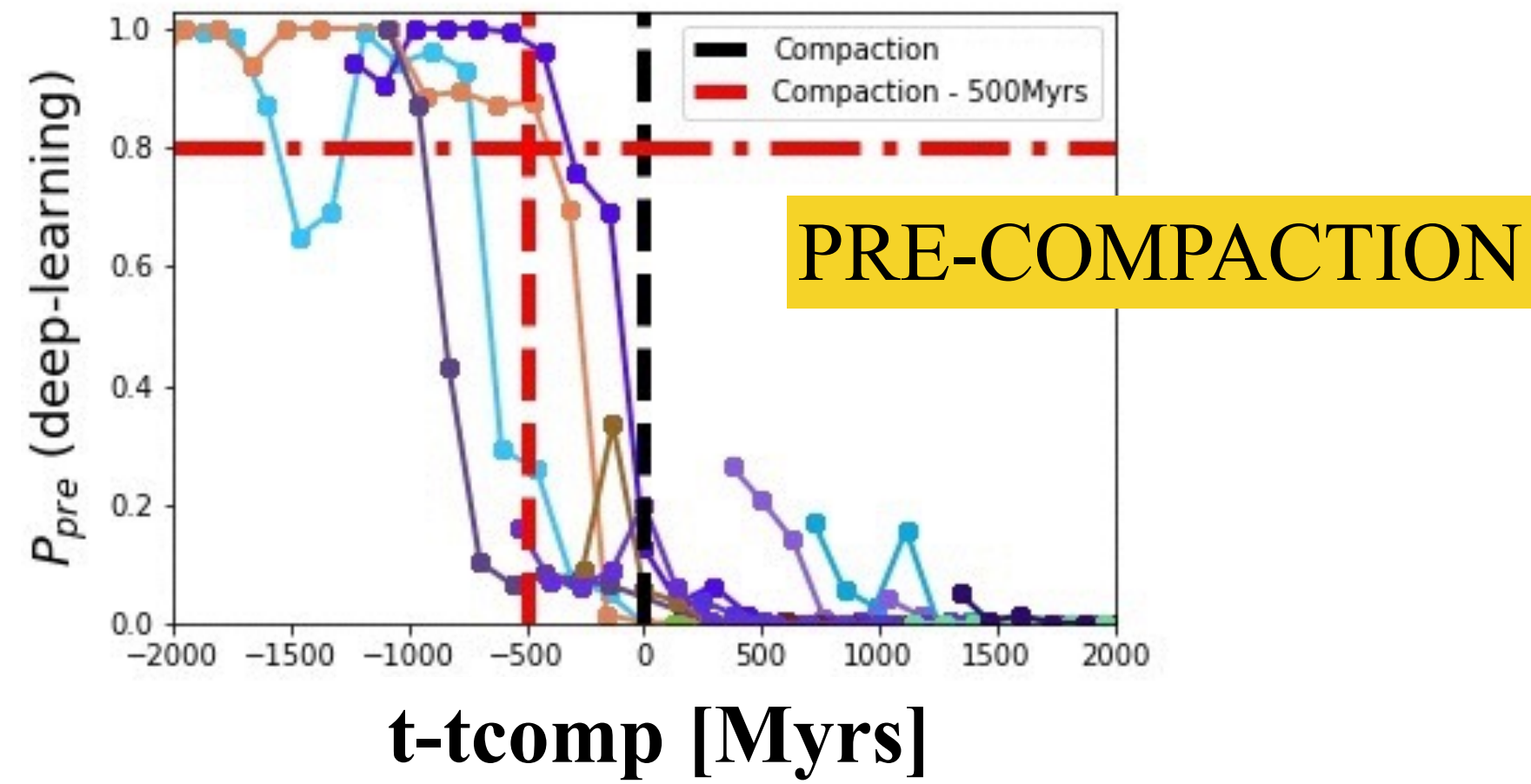
**MOCK  
 HST  
 IMAGES  
 (CANDELS  
 filters)**



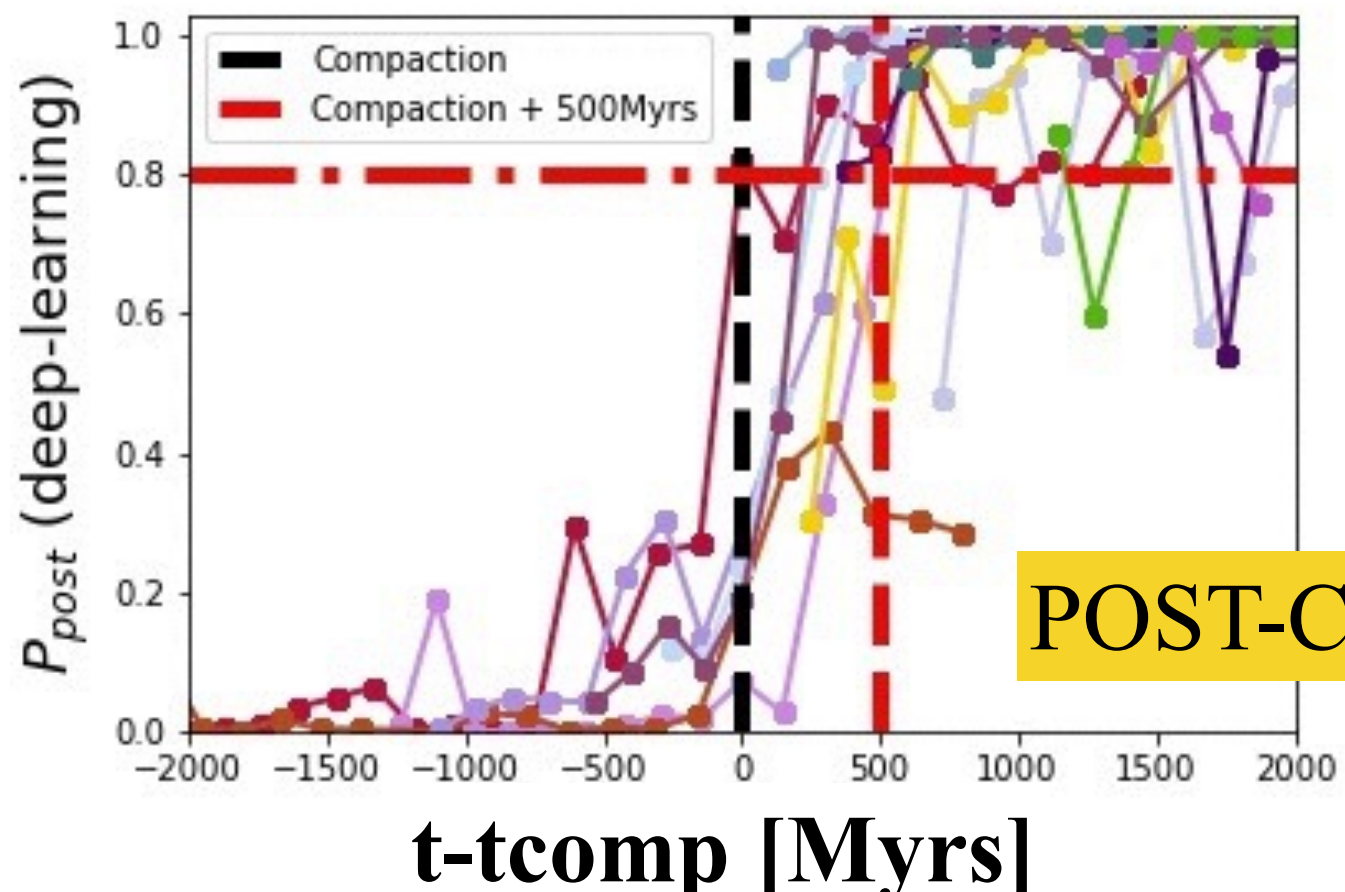
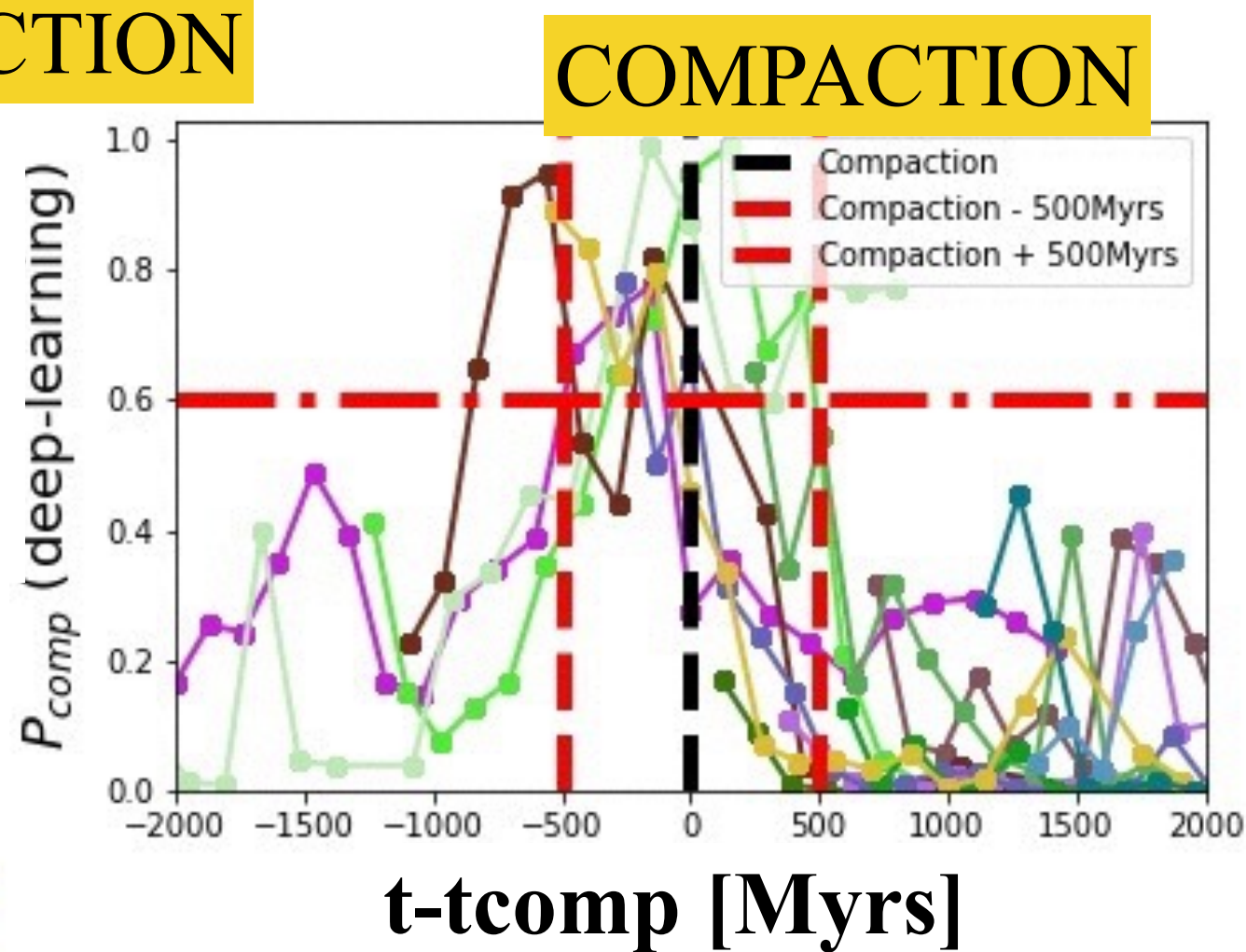
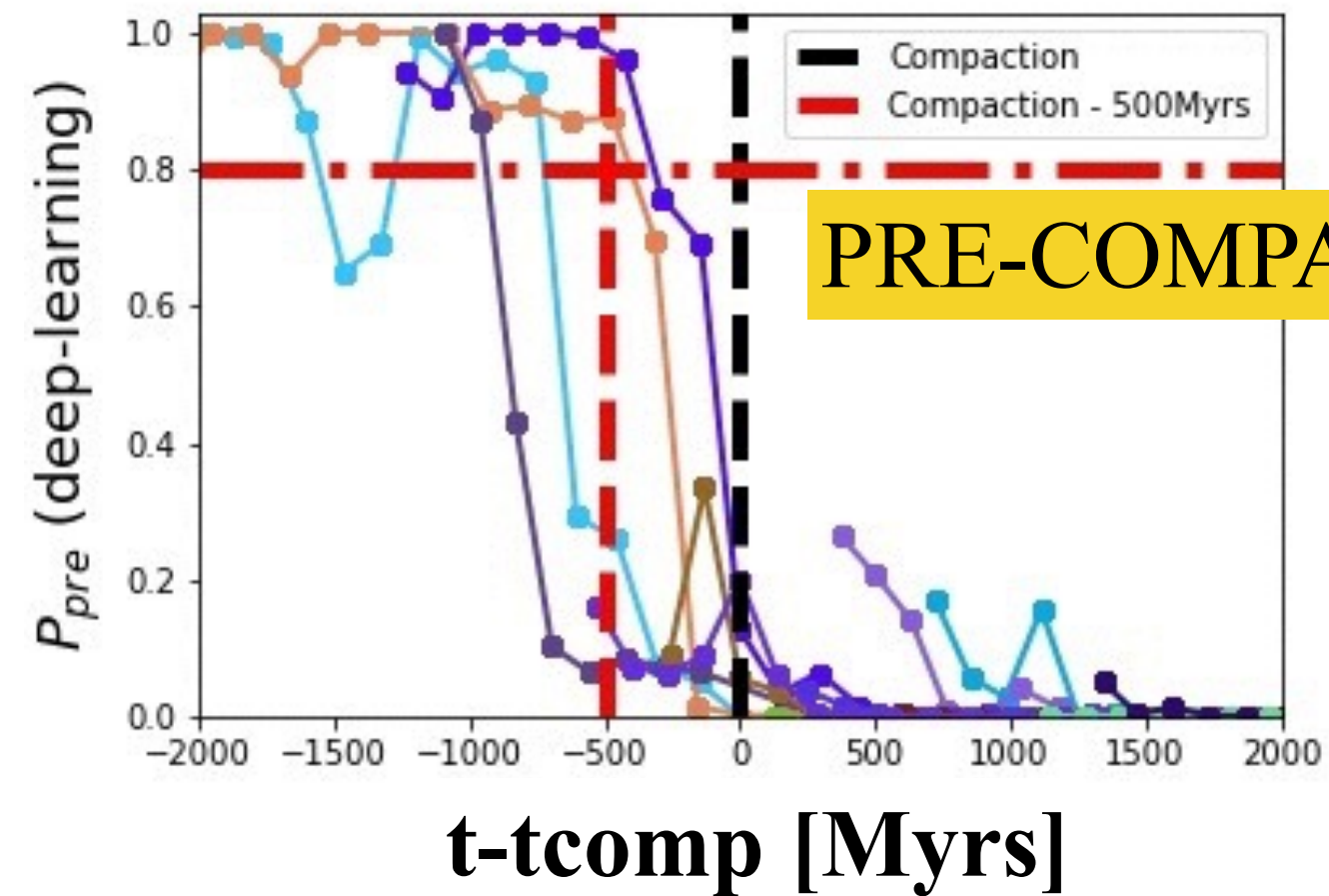
**Pre-Compaction**

**Compaction**

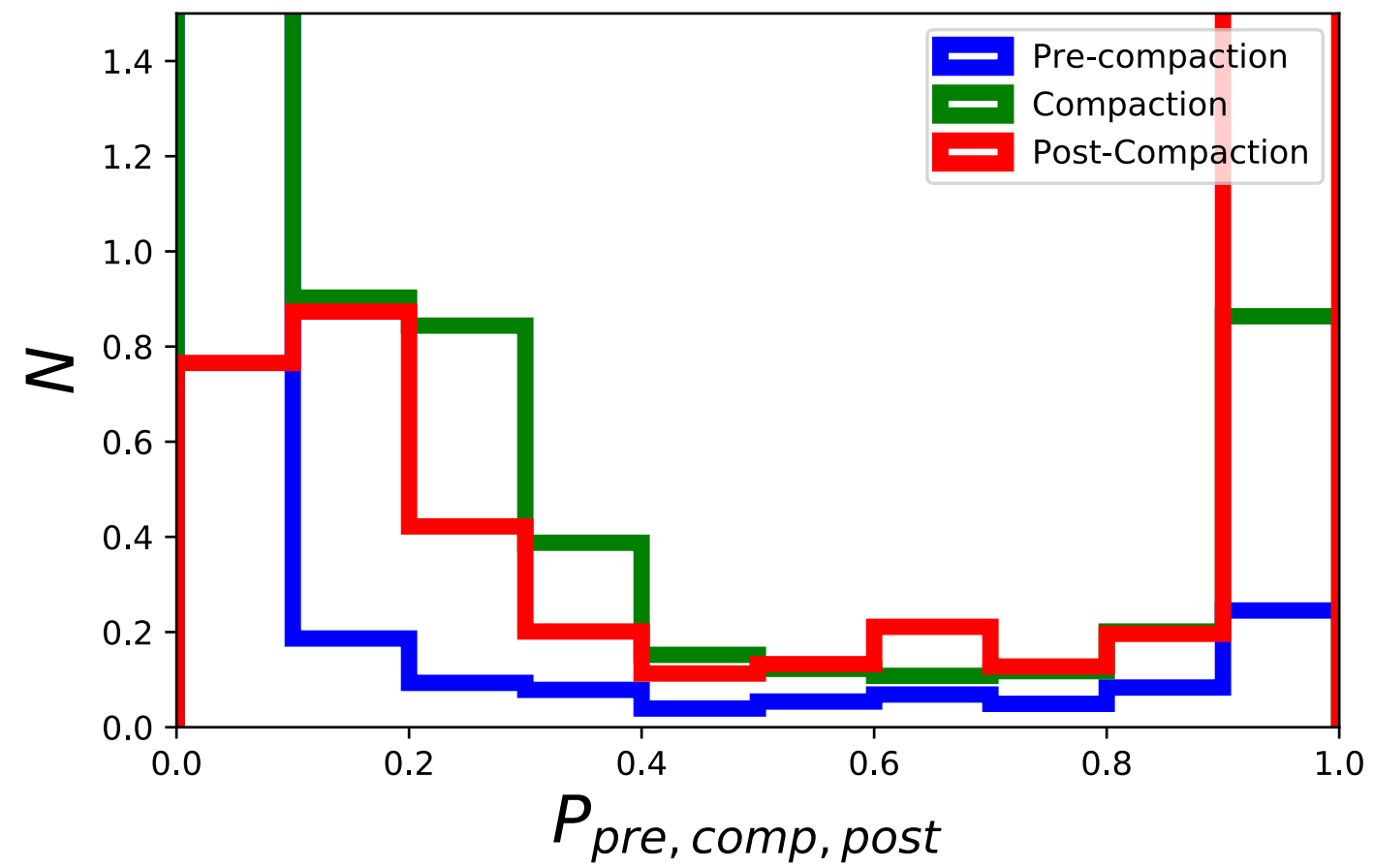
**Post-Compaction**





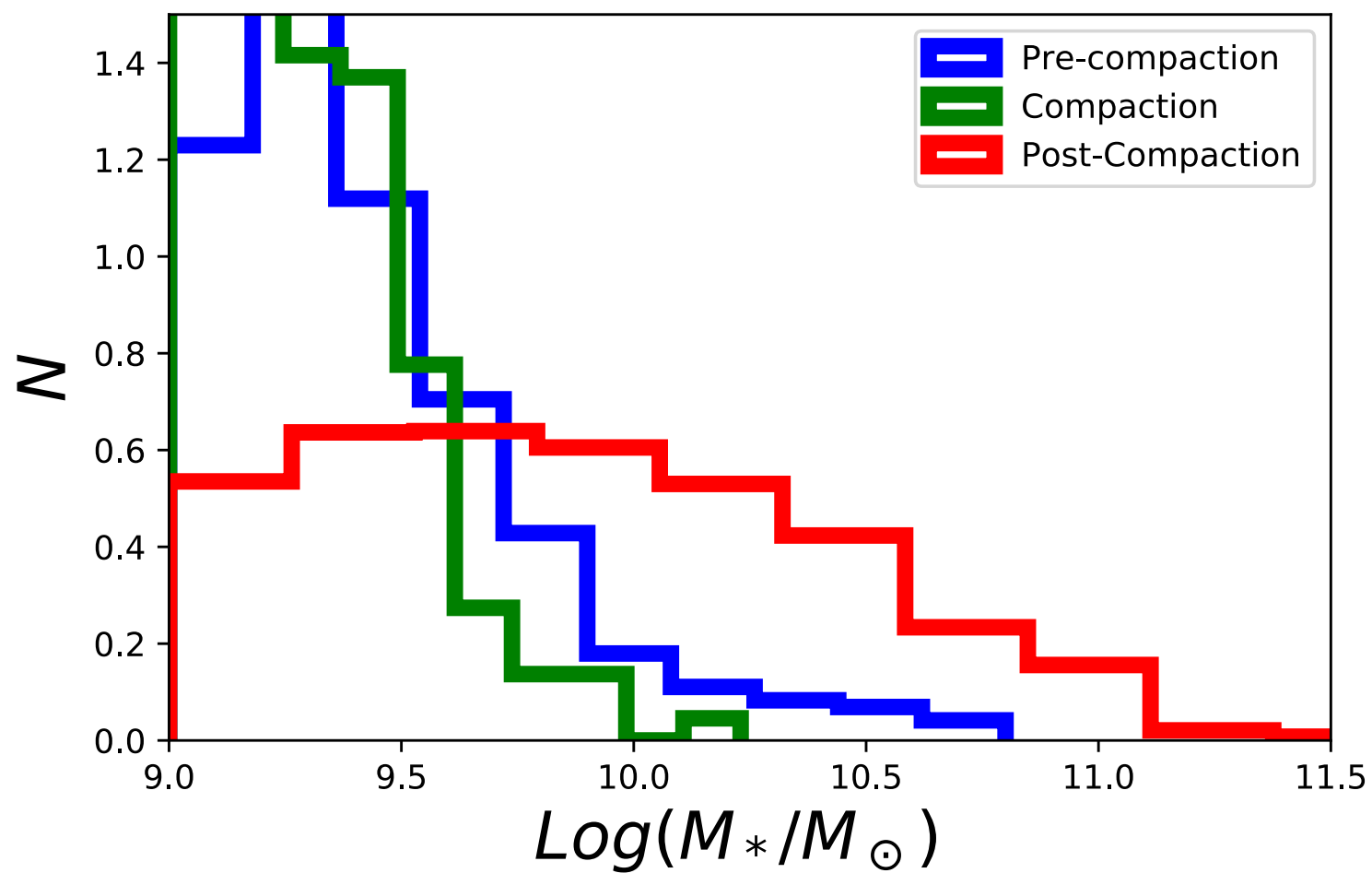
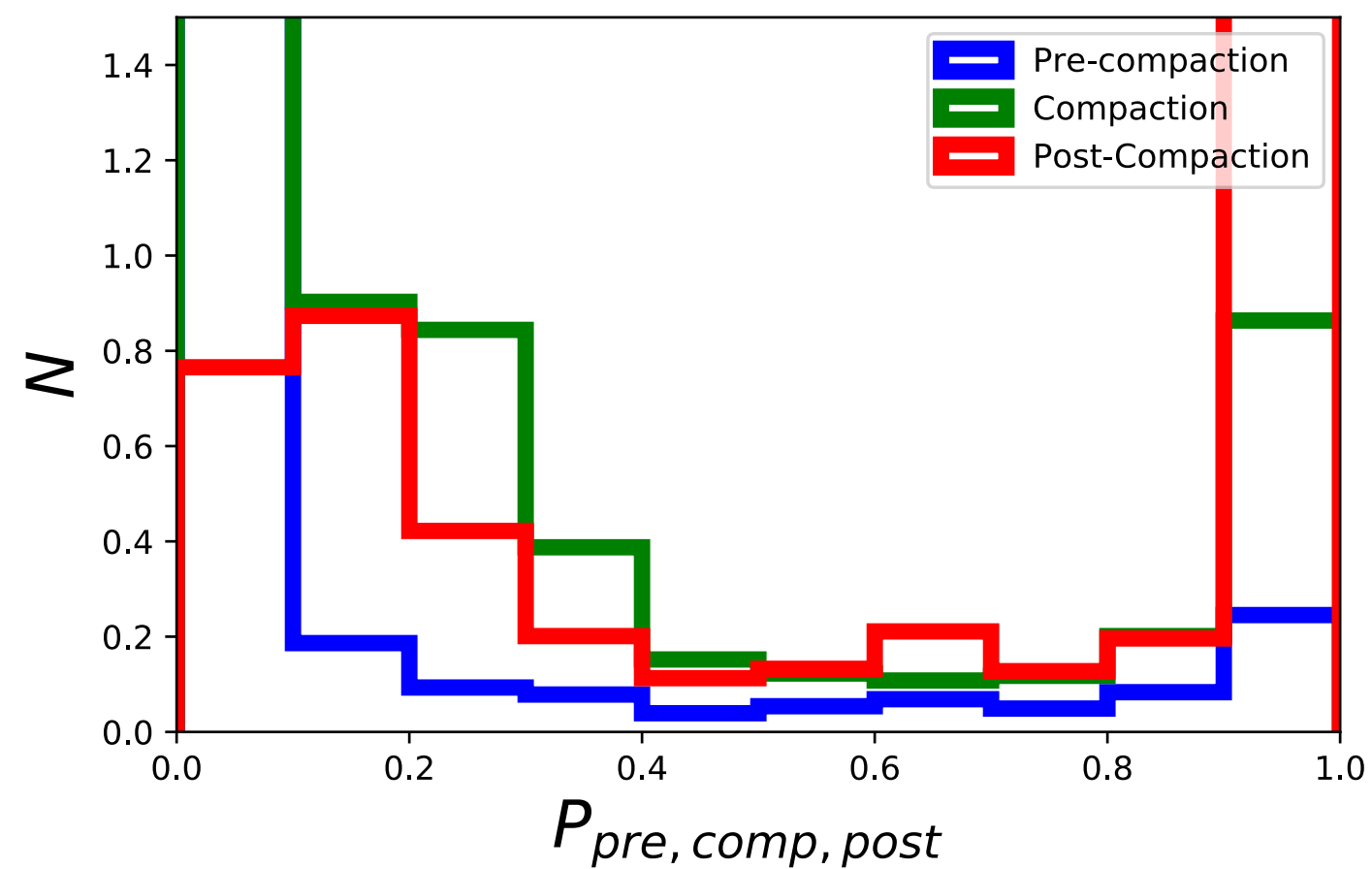


# WHEN APPLIED TO REAL CANDELS DATA

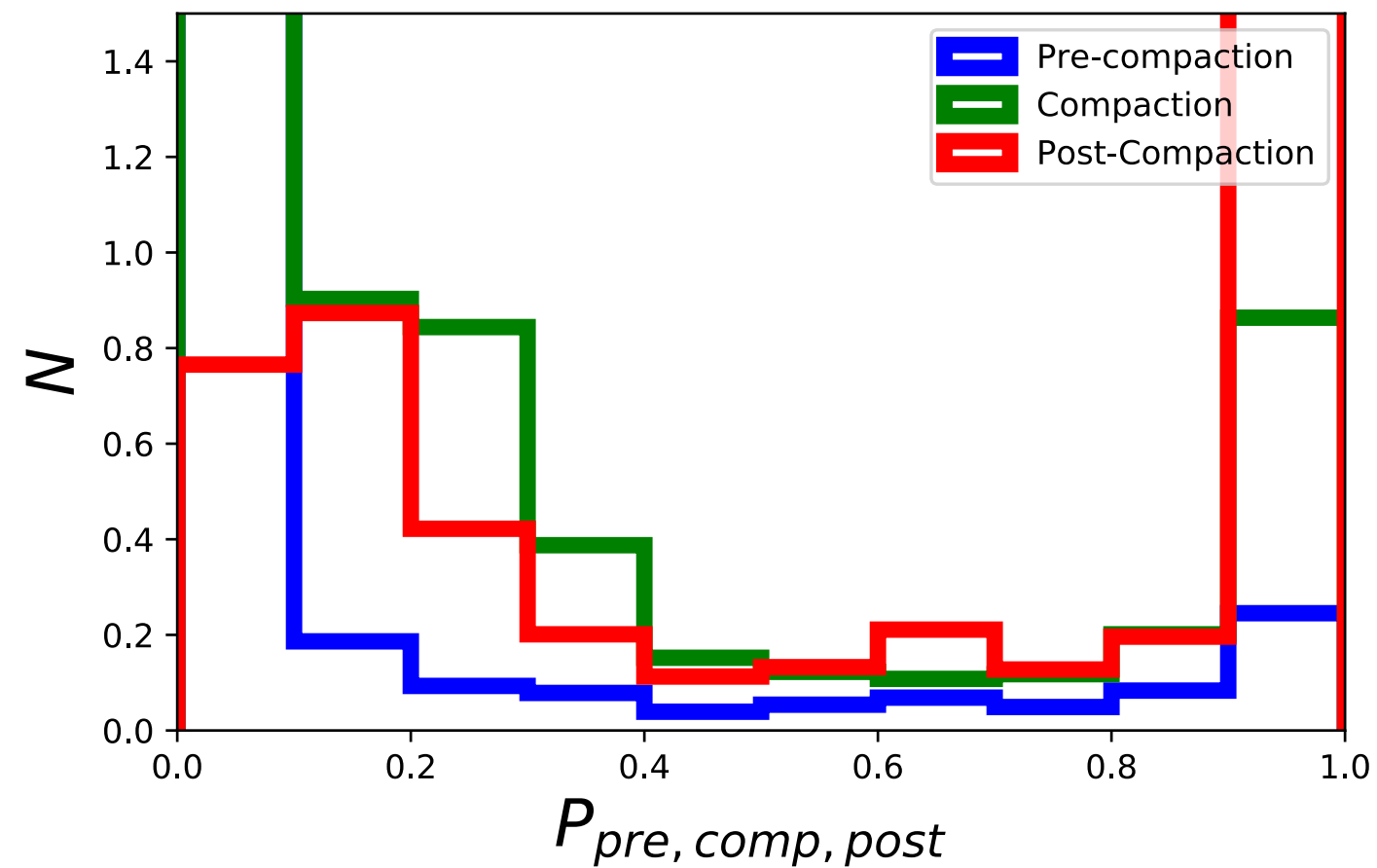




# WHEN APPLIED TO REAL CANDELS DATA

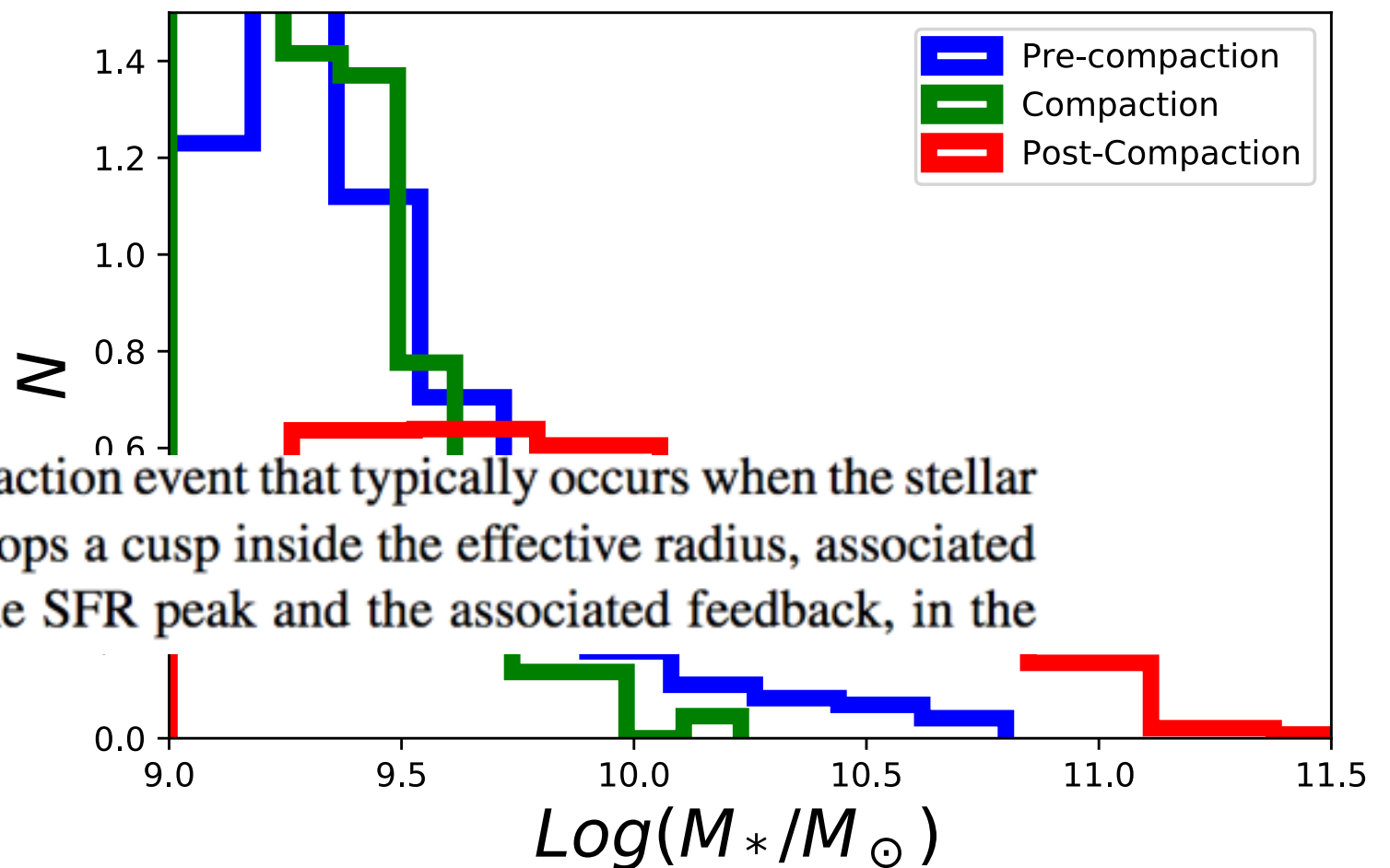


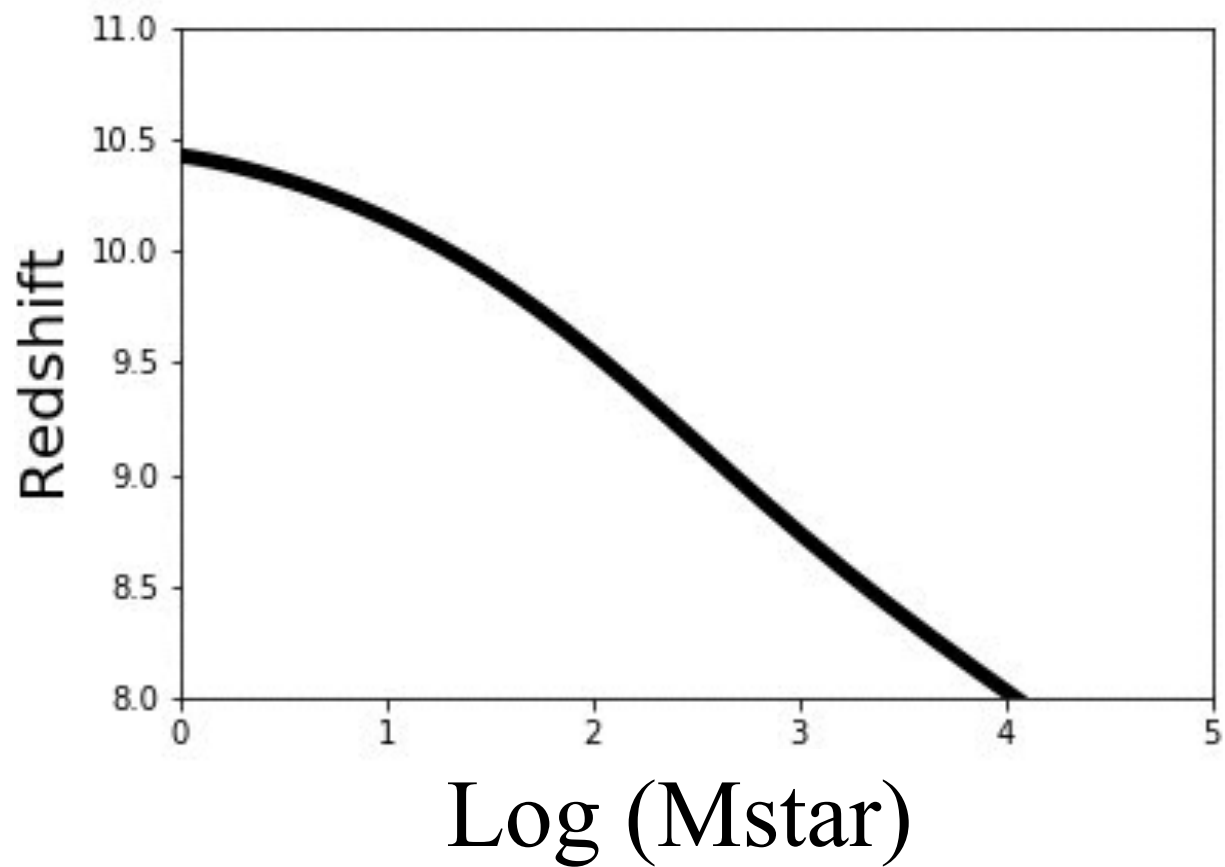
# WHEN APPLIED TO REAL CANDELS DATA



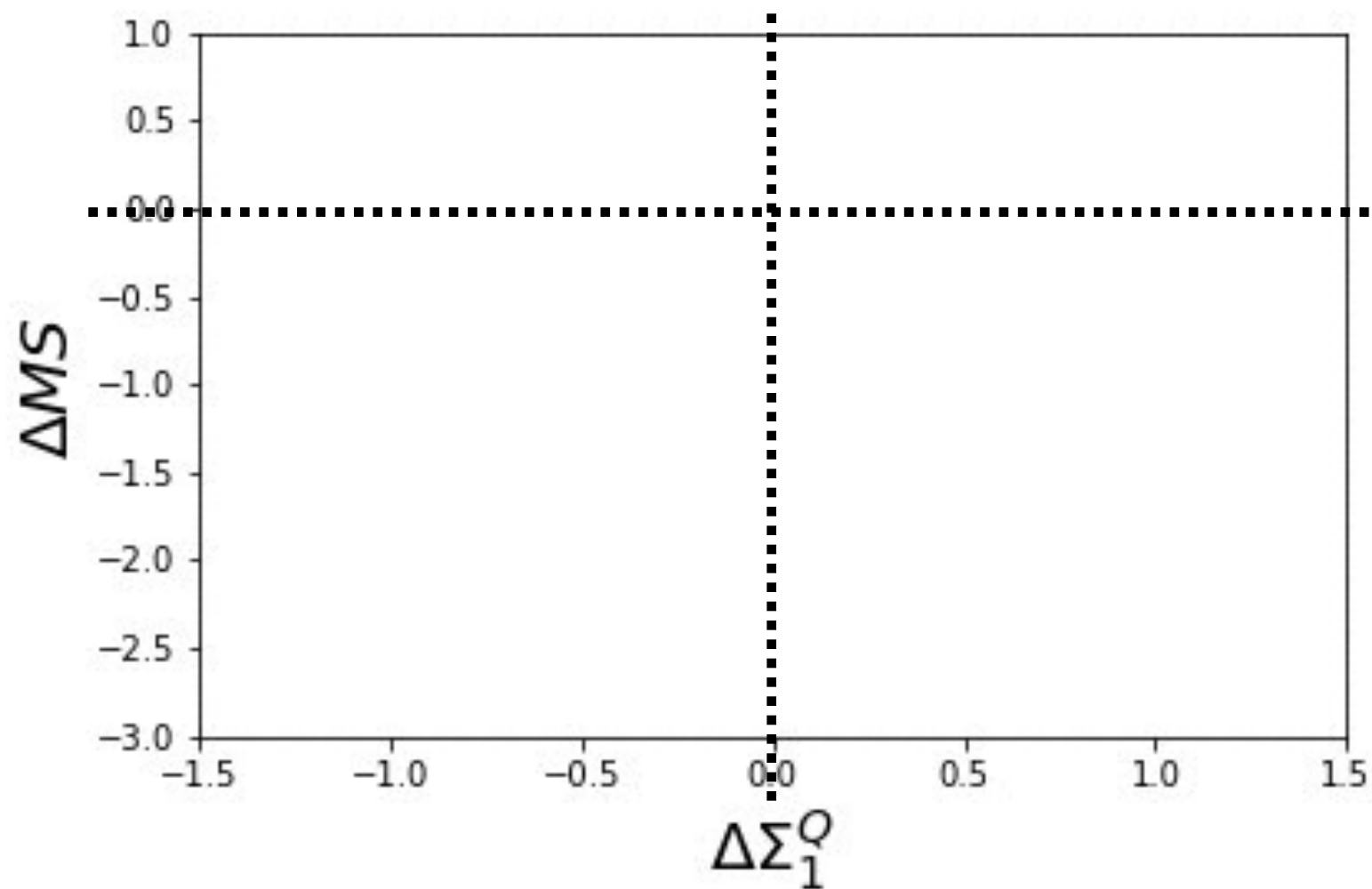
their evolution phases. Following a wet compaction event that typically occurs when the stellar mass is  $\sim 10^{9.5} M_{\odot}$  at  $z \sim 2-4$ , the gas develops a cusp inside the effective radius, associated with a peak in star formation rate (SFR). The SFR peak and the associated feedback, in the

**Tacchella+16**

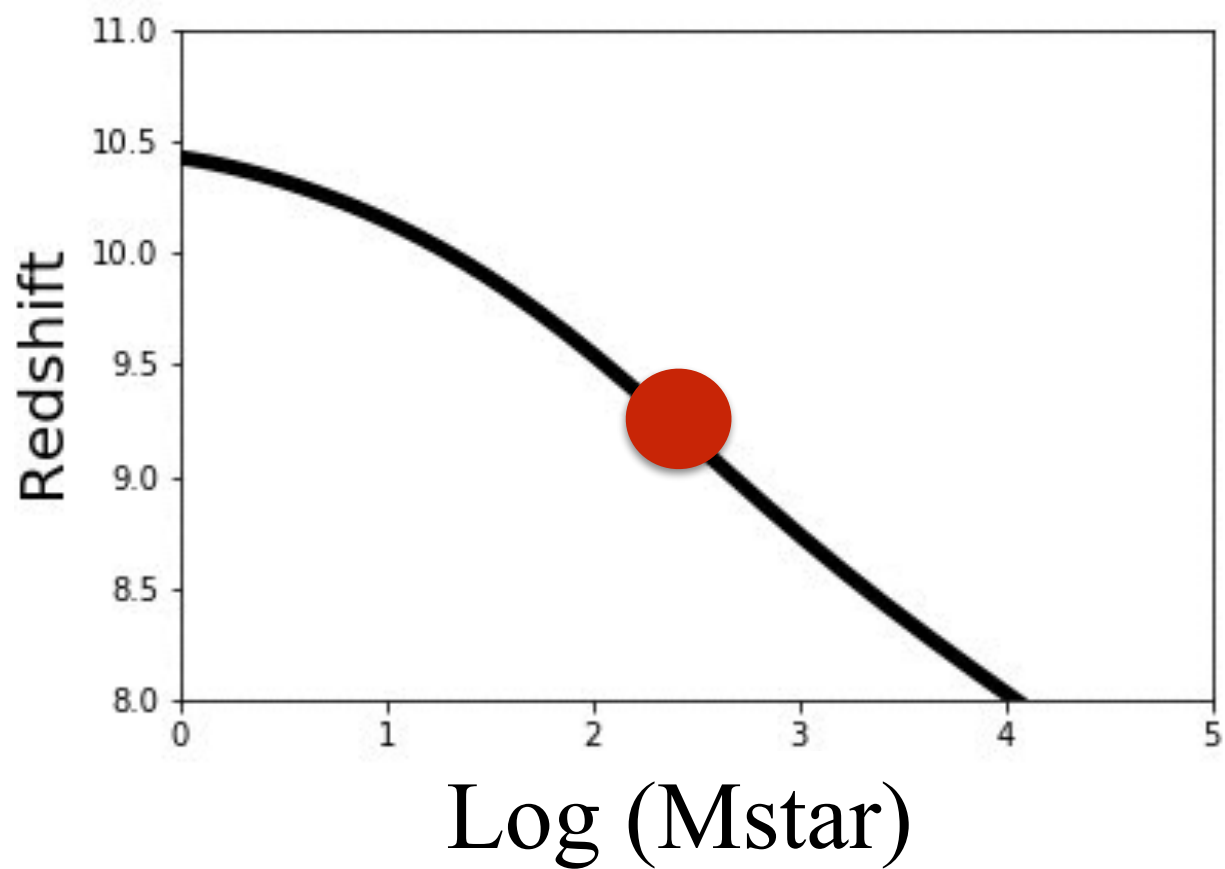




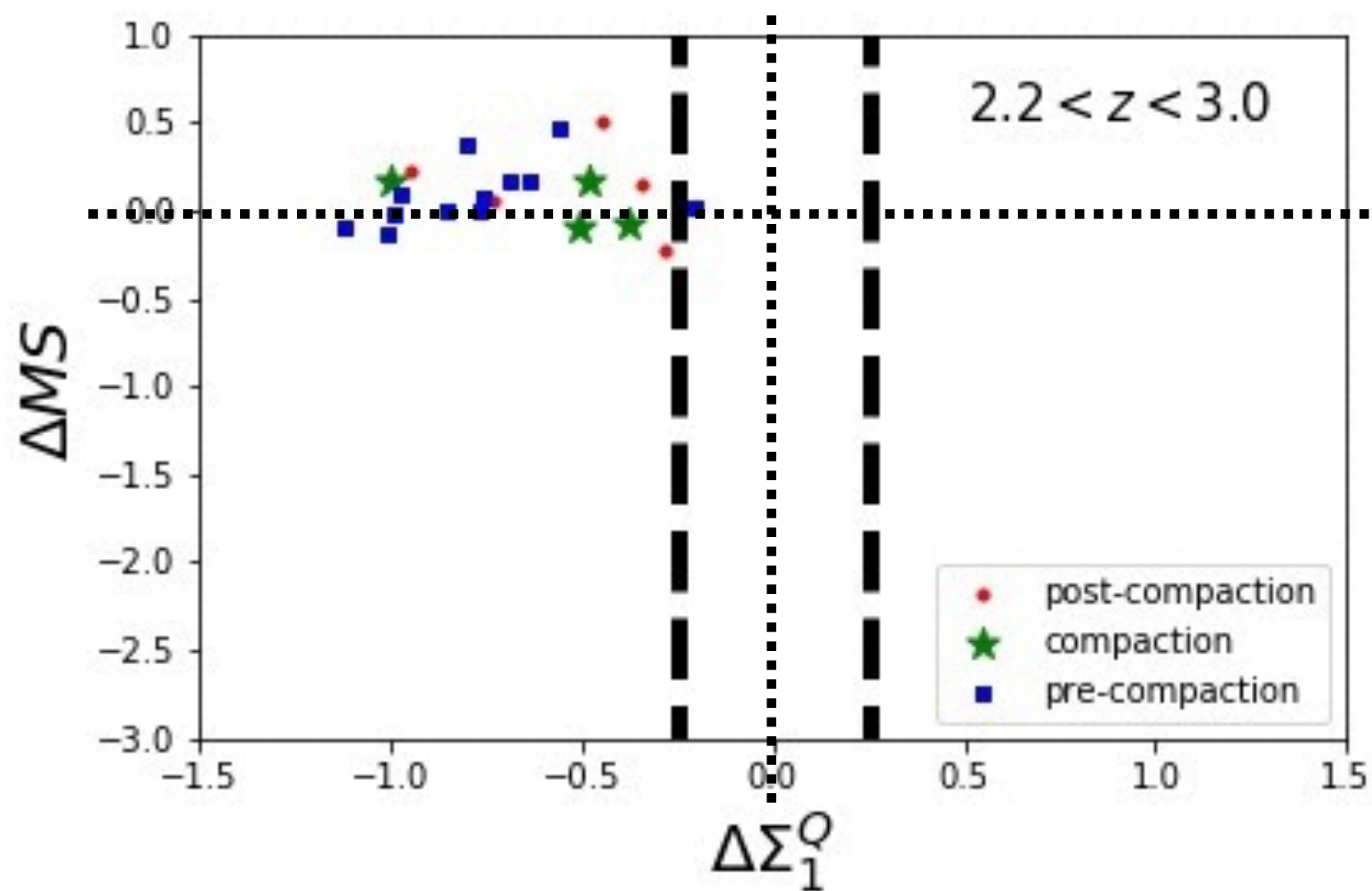
**abundance matching  
[Rodriguez-Puebla+17]**



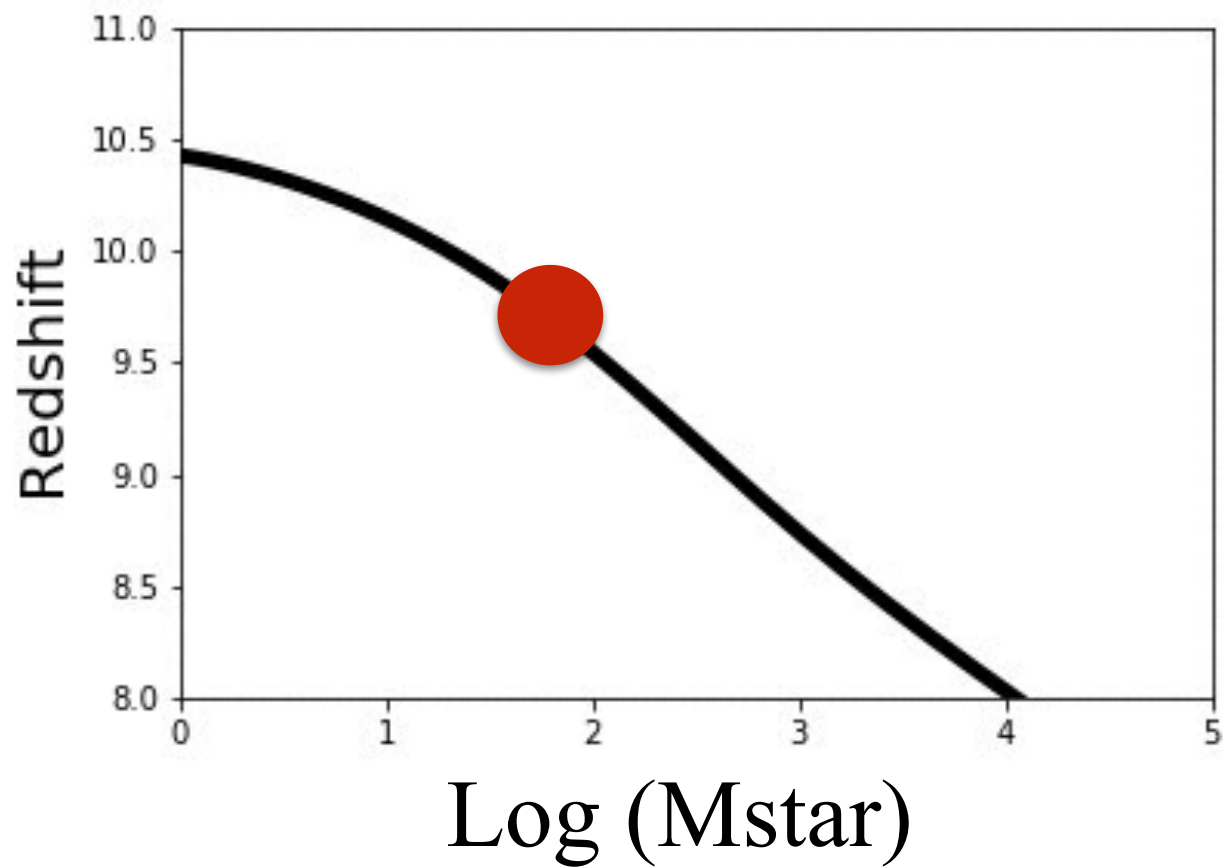
**MHC+18 (in prep)**



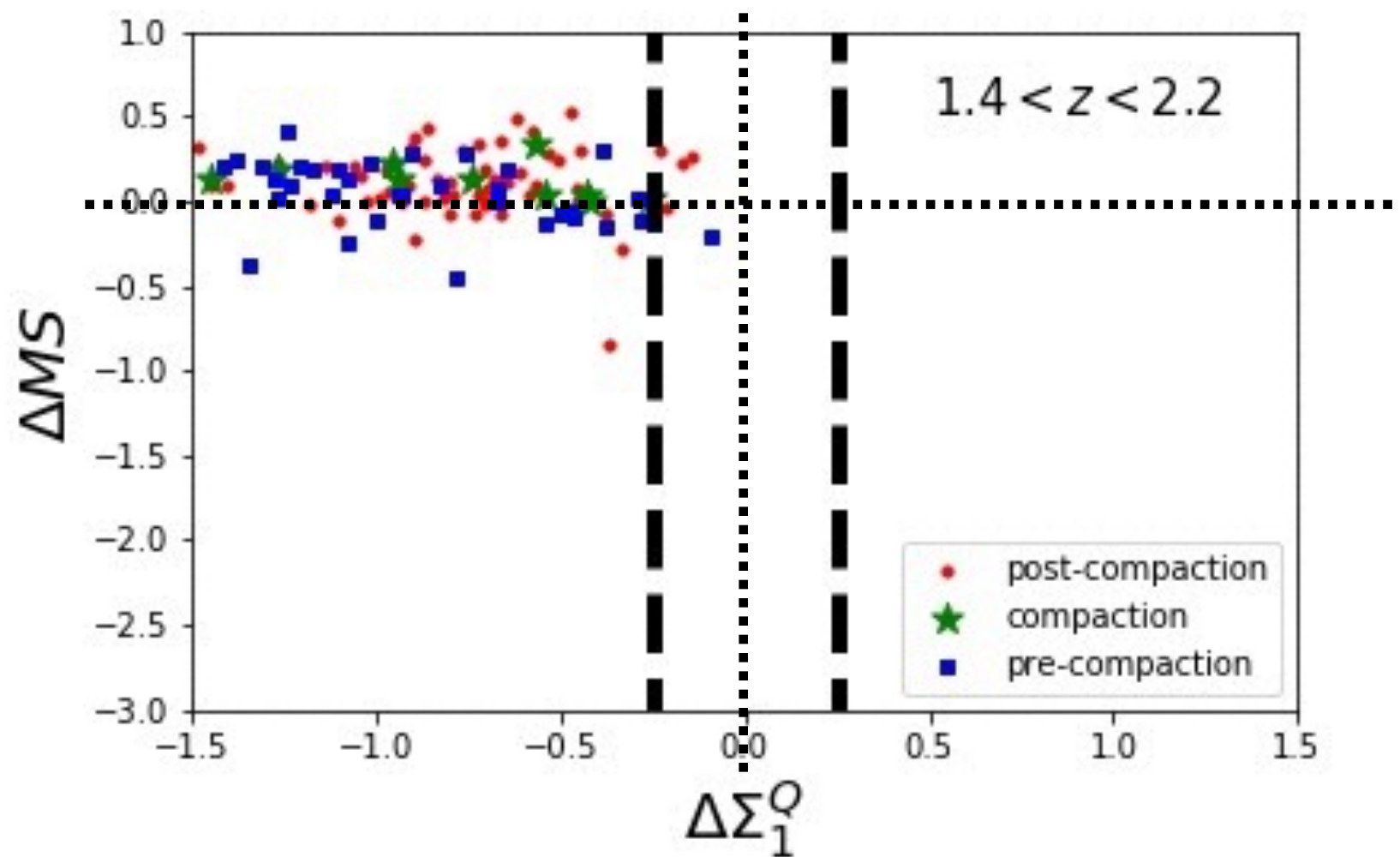
**abundance matching  
[Rodriguez-Puebla+17]**



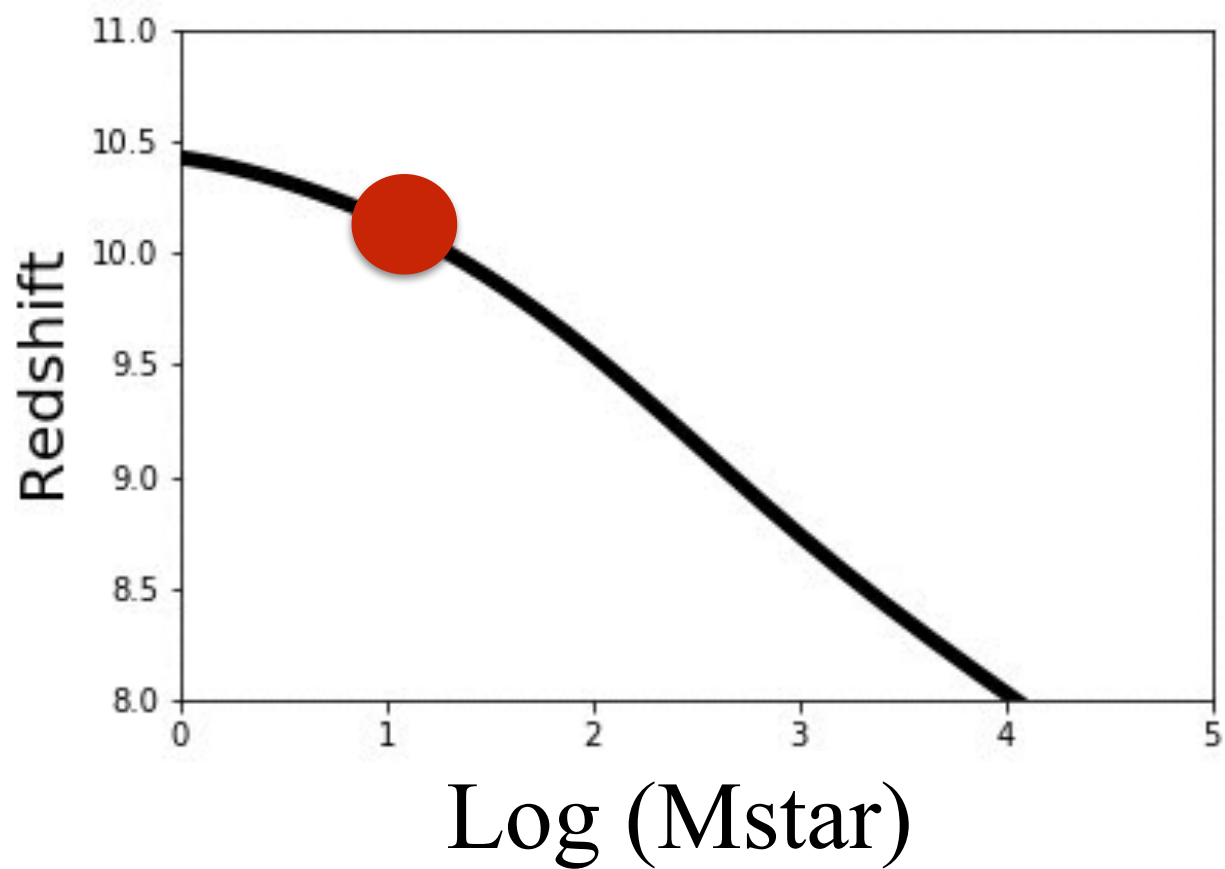
**MHC+18 (in prep)**



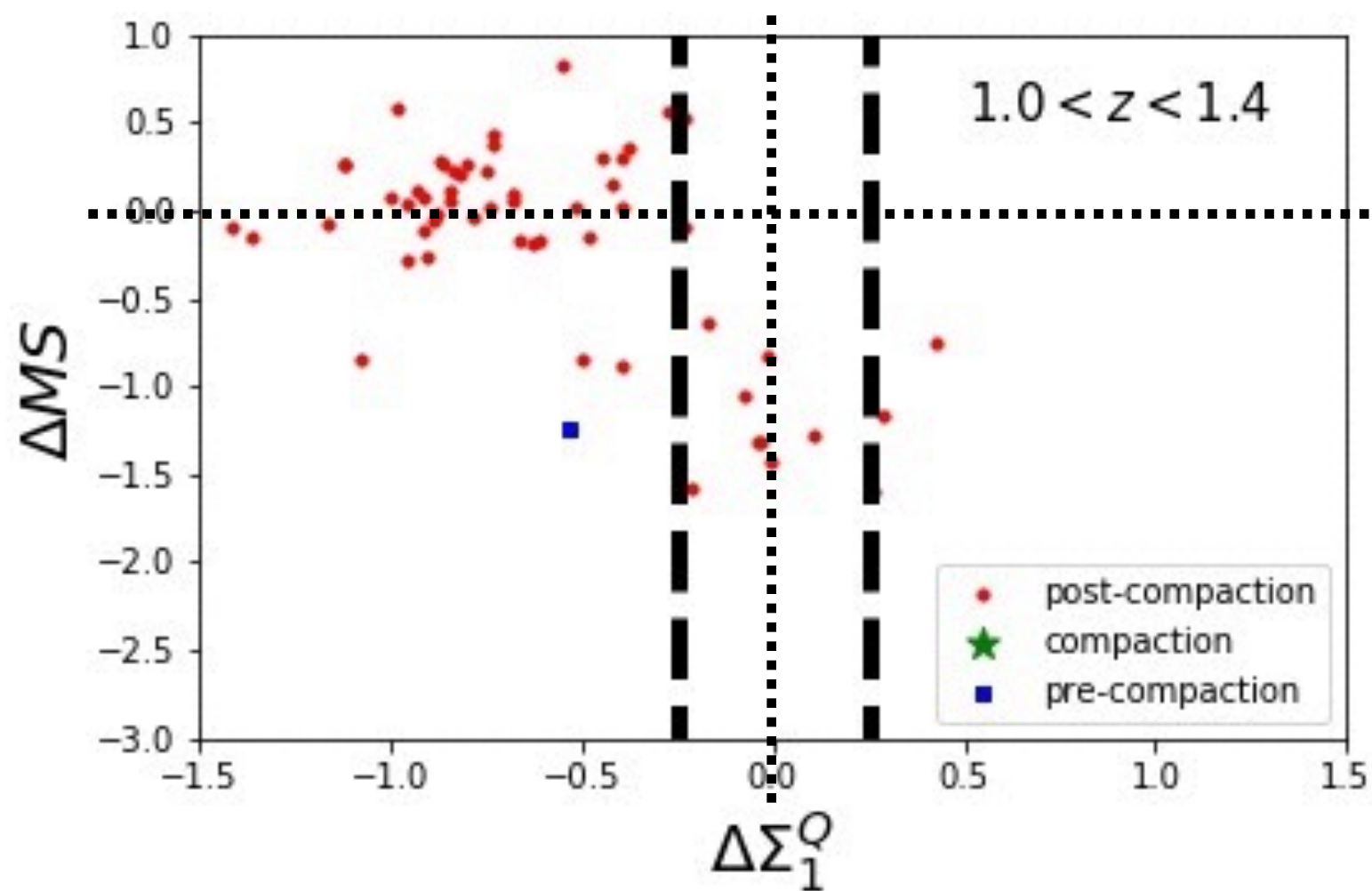
**abundance matching  
[Rodriguez-Puebla+17]**



**MHC+18 (in prep)**

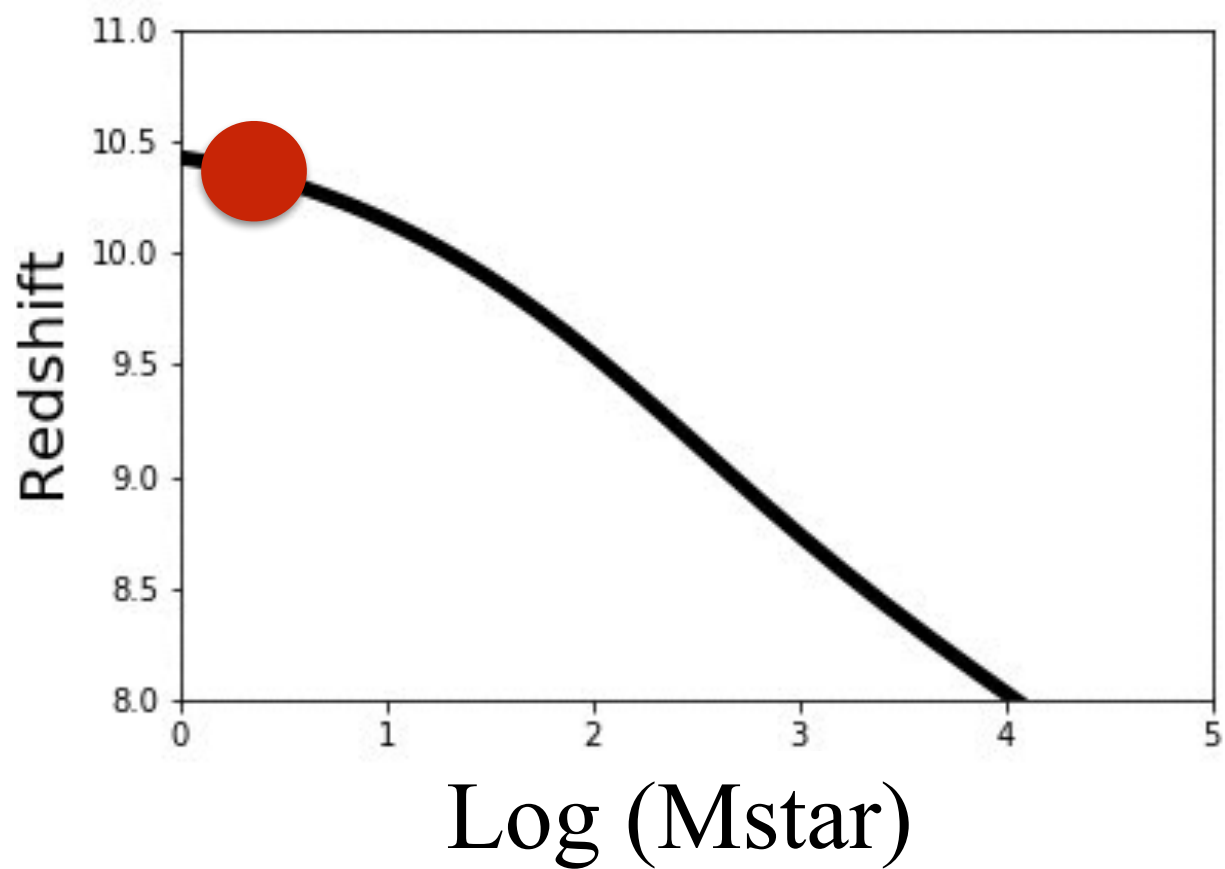


**abundance matching  
[Rodriguez-Puebla+17]**

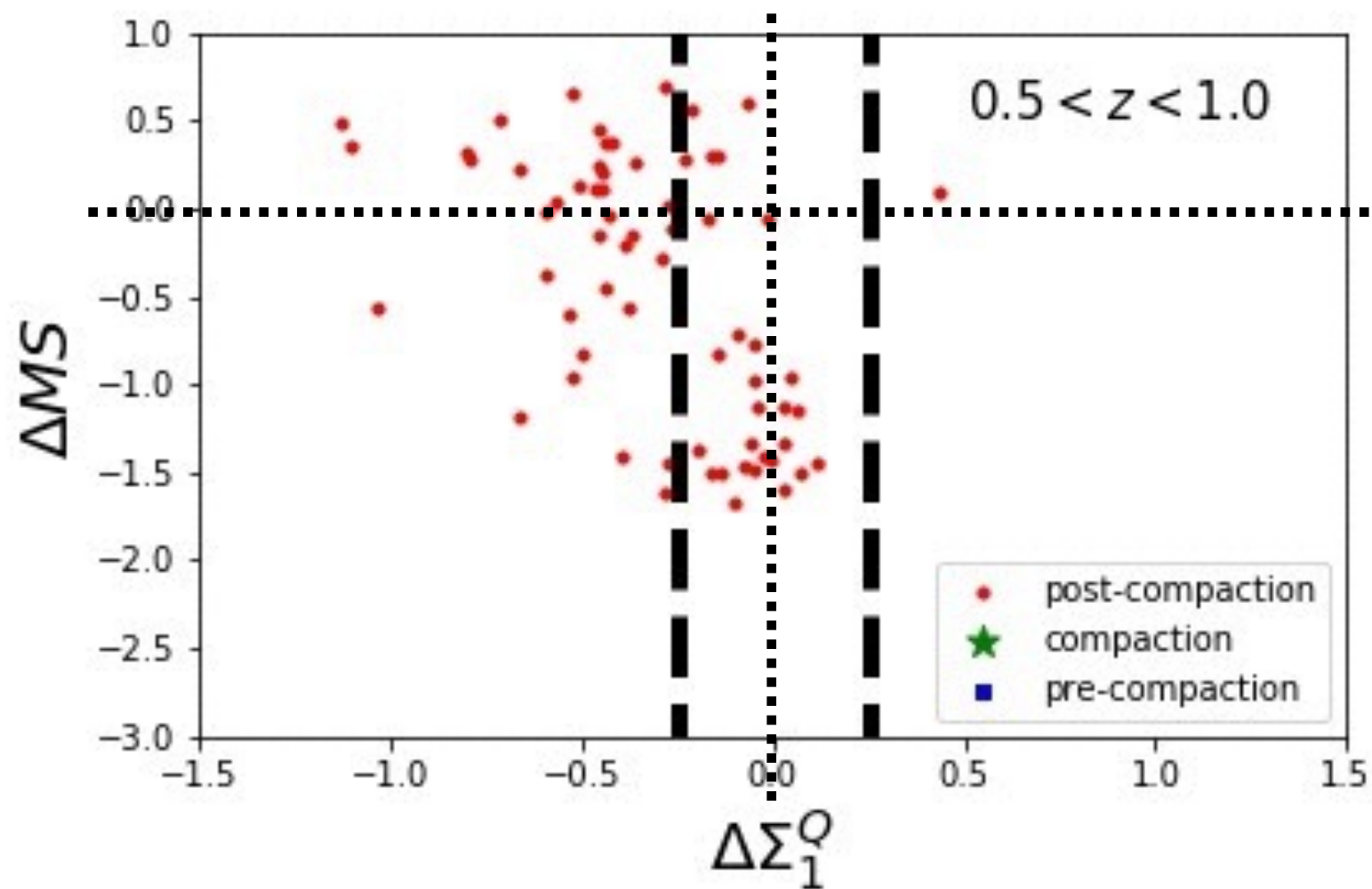


**MHC+18 (in prep)**



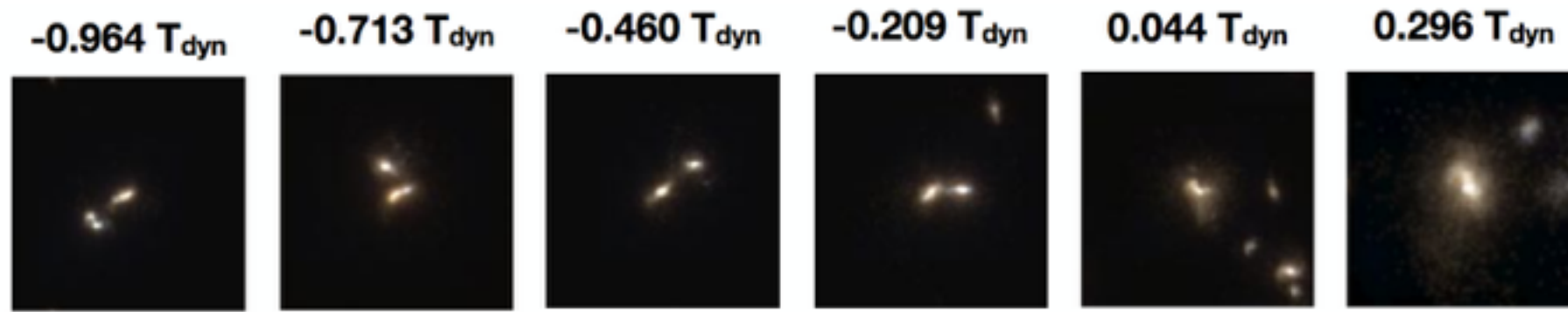


**abundance matching  
[Rodriguez-Puebla+17]**

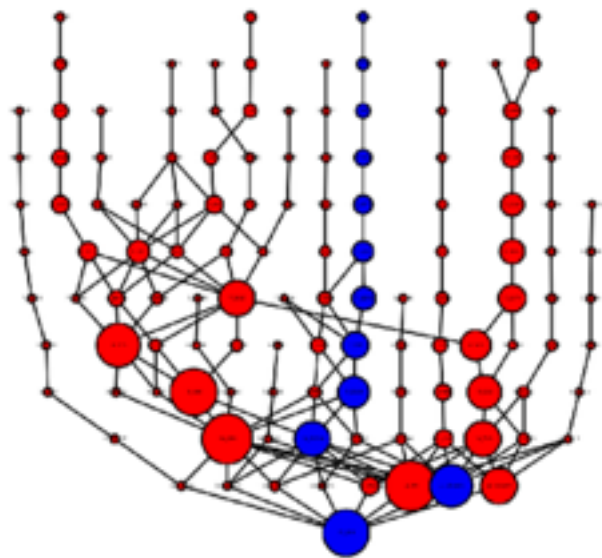
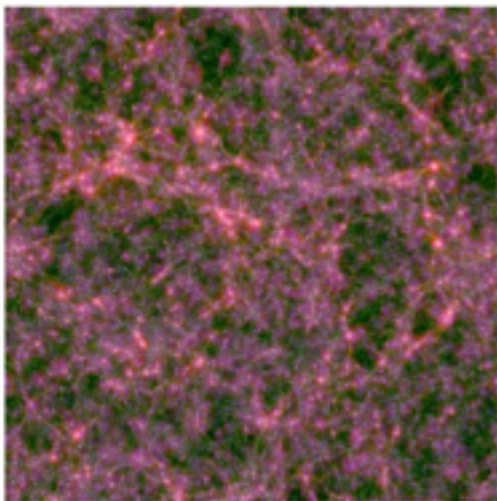


**MHC+18 (in prep)**

# Dissecting mergers with DL



**Horizon AGN  
hydro sim**  
[Dubois+14]

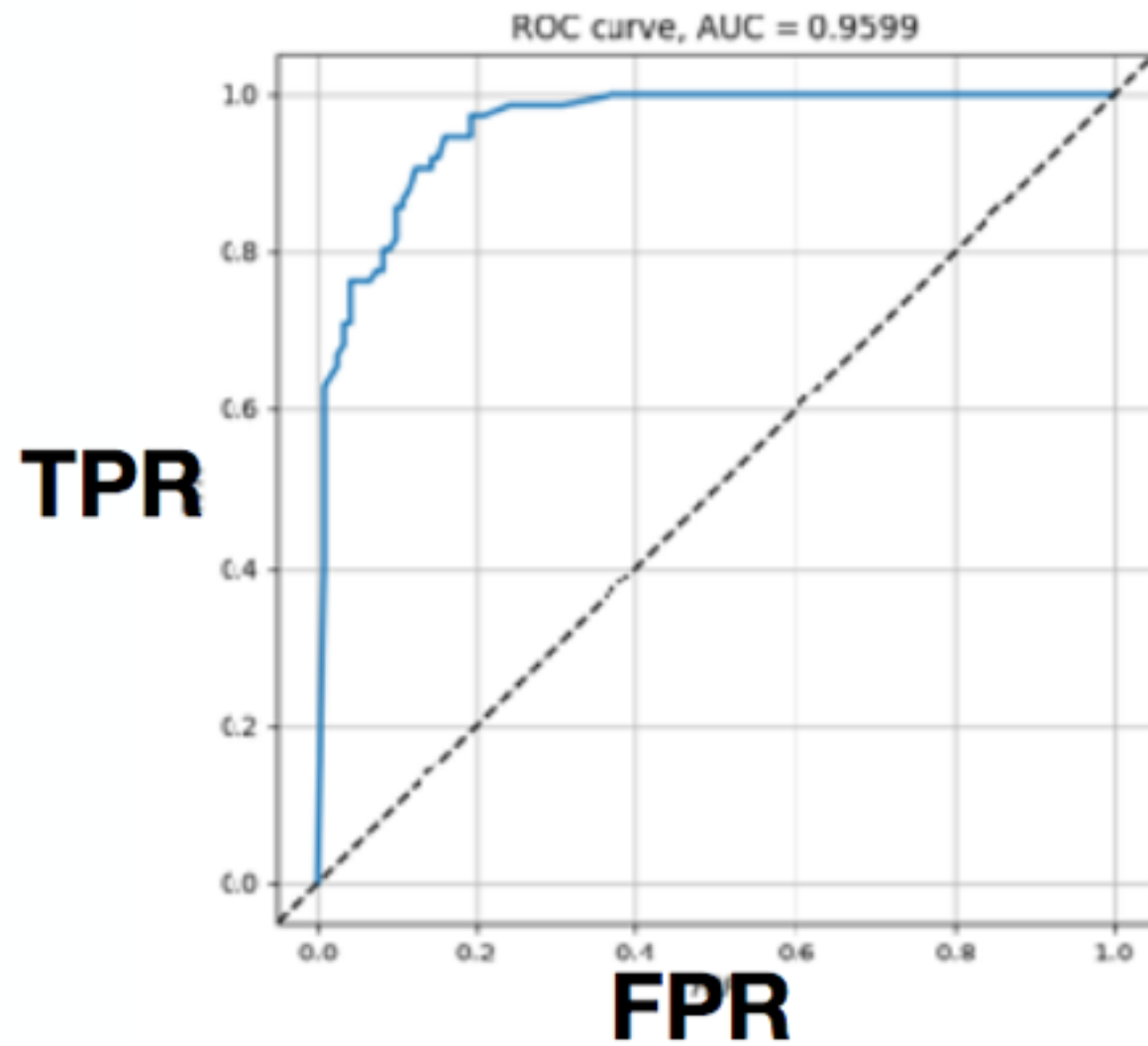


Mock images

Deep  
Learning

Merger properties:  
- mass-ratio  
- stage  
- timescale

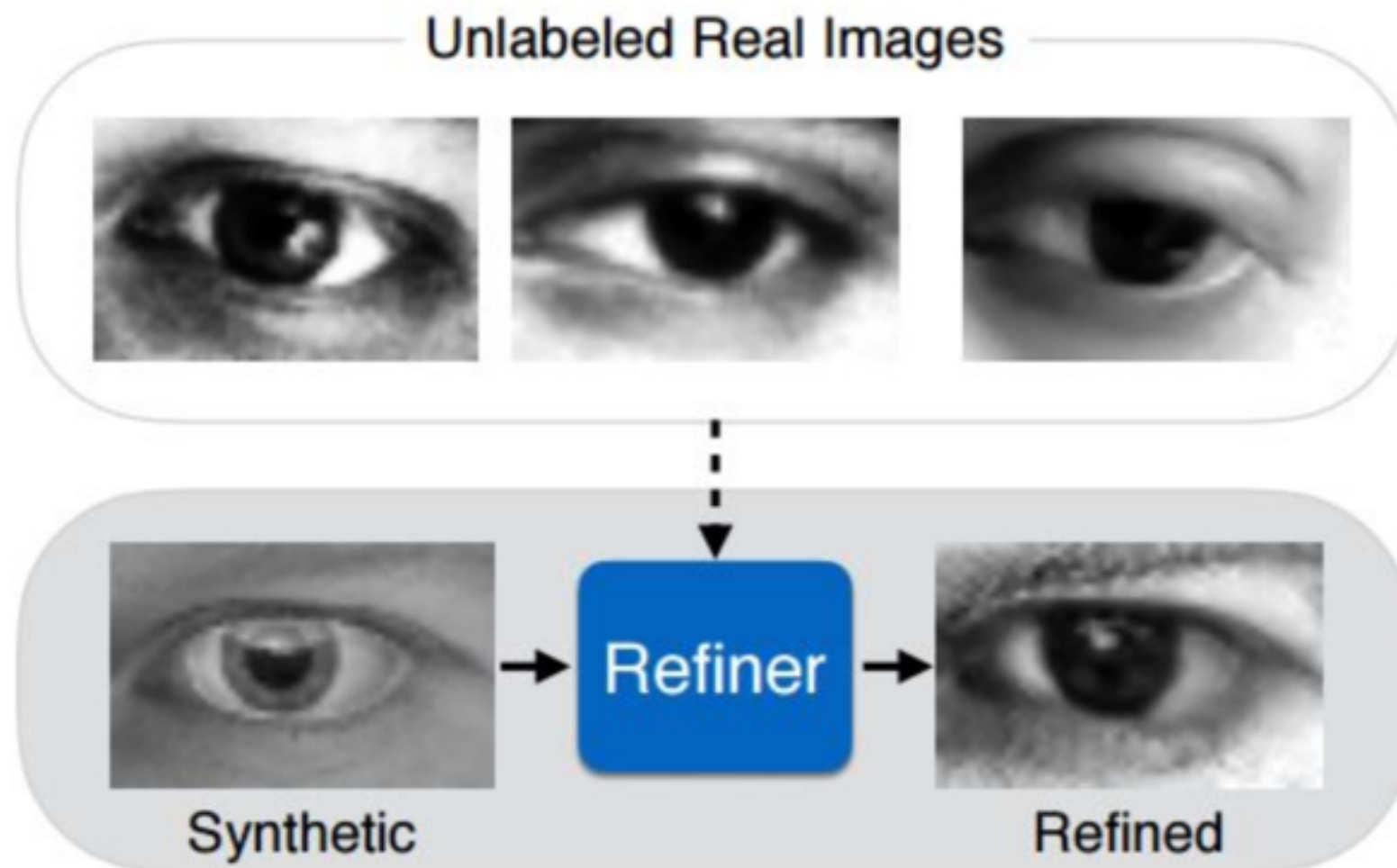
# MERGER PHASE



**Accuracy over validation set: ~ 96%**

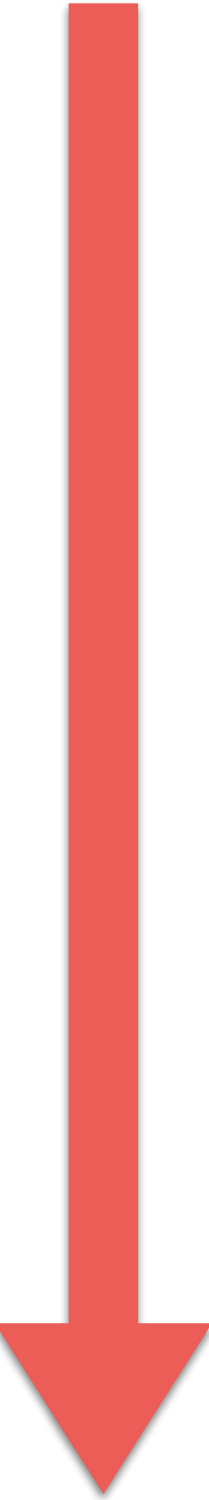
- **GROUP #4:** Finding the unknown?

[ANR project submitted]



# DL FOR GALAXIES?

- **GROUP #1:** Time consuming tasks that humans do easily but classically challenging for computers - classification of objects
- **GROUP #2:** Efficient and fast quantitative measurements on large amount of (multi-lambda) data [photoz's, sizes, ellipticities]
- **GROUP #3:** Find hidden unknown correlations in the data, - Linking observations and theory
- **GROUP #4:** Finding the unknown?

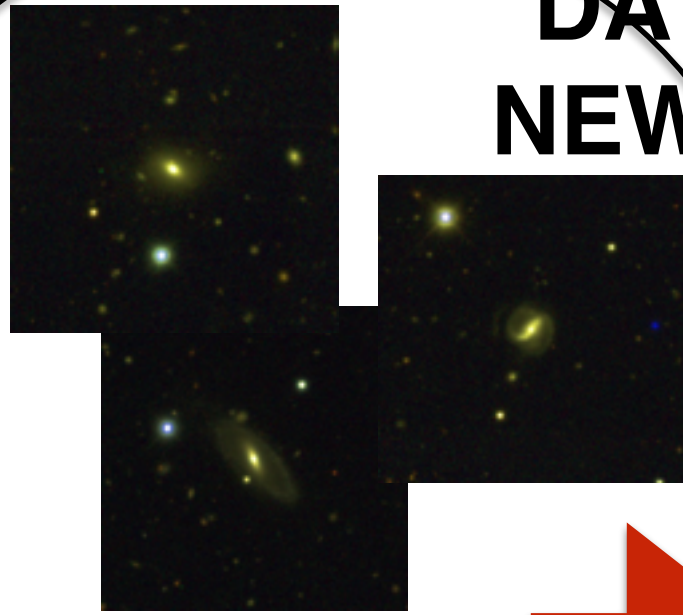


# Group #1: Classification of large datasets



- **Requires a huge volume** of “labeled” data to train.  
Human intervention is necessary anyway...

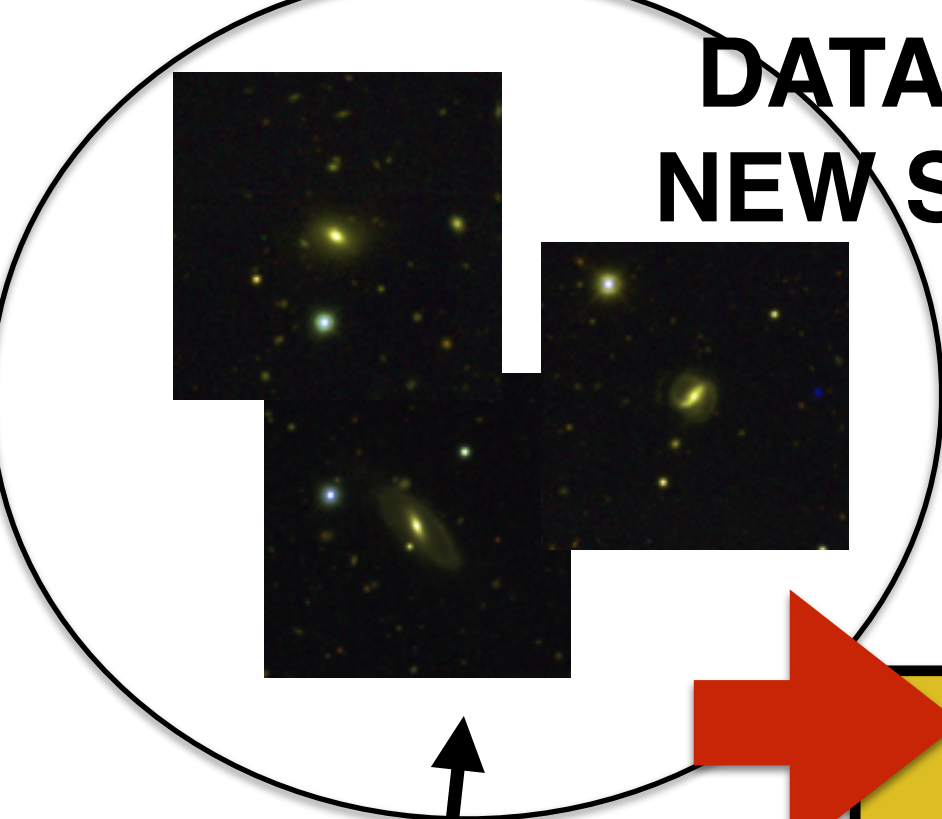
**DATA FROM  
NEW SURVEY**



**Requires a huge volume** of “labeled” data to train. Human intervention is necessary anyway.

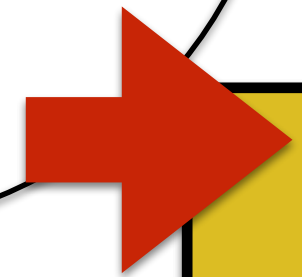
DEEP-LEARNING  
BASED  
MACHINE

Galaxy ZOO like  
classifications for  
for the entire  
sample

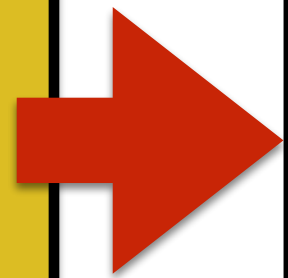


**DATA FROM  
NEW SURVEY**

**Requires a huge volume** of “labeled” data to train. Human intervention is necessary anyway.



DEEP-LEARNING  
BASED  
MACHINE



Galaxy ZOO like classifications for for the entire sample

**Human classifications**



**DATA FROM  
NEW SURVEY**

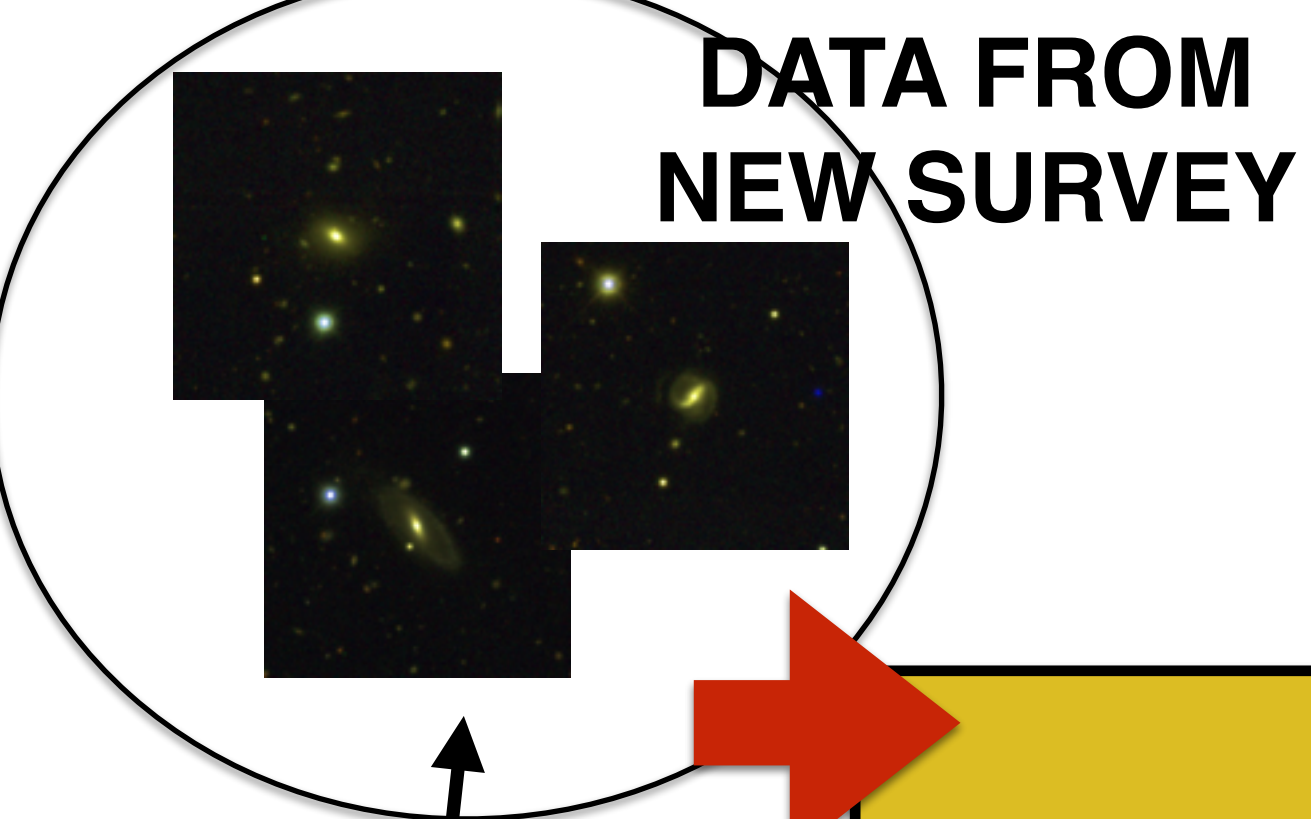
**Requires a huge volume** of  
“labeled” data to train.  
Human intervention is  
necessary anyway.

**N??**

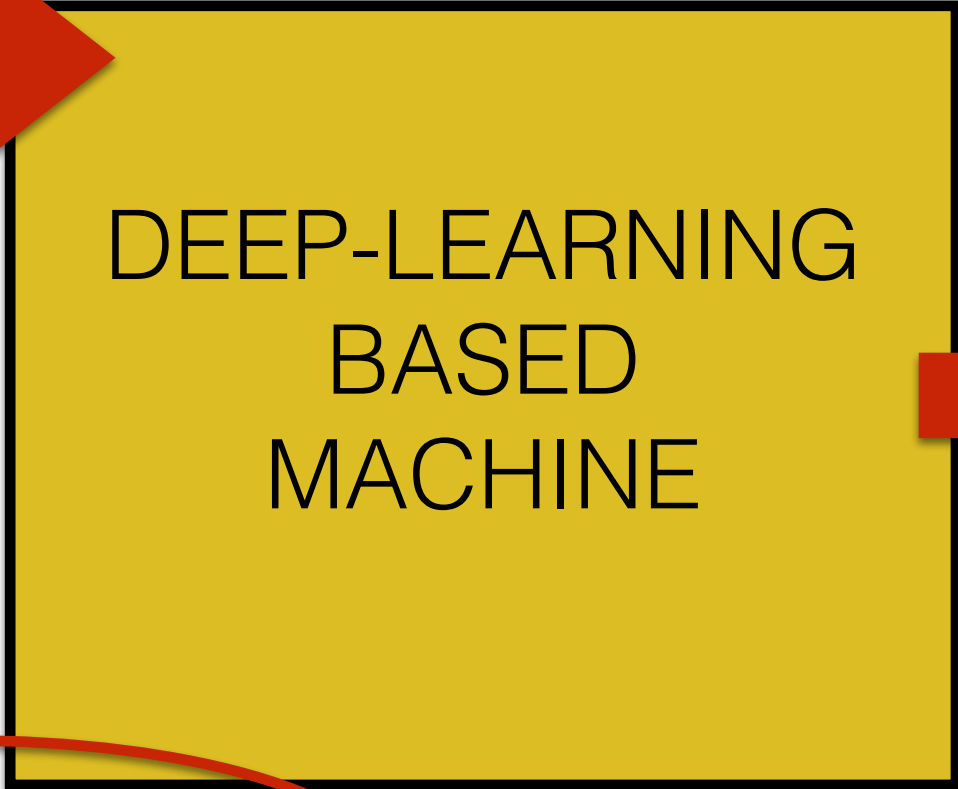
DEEP-LEARNING  
BASED  
MACHINE

Galaxy ZOO like  
classifications for  
for the entire  
sample

**Human classifications**



**Requires a huge volume** of “labeled” data to train. Human intervention is necessary anyway.

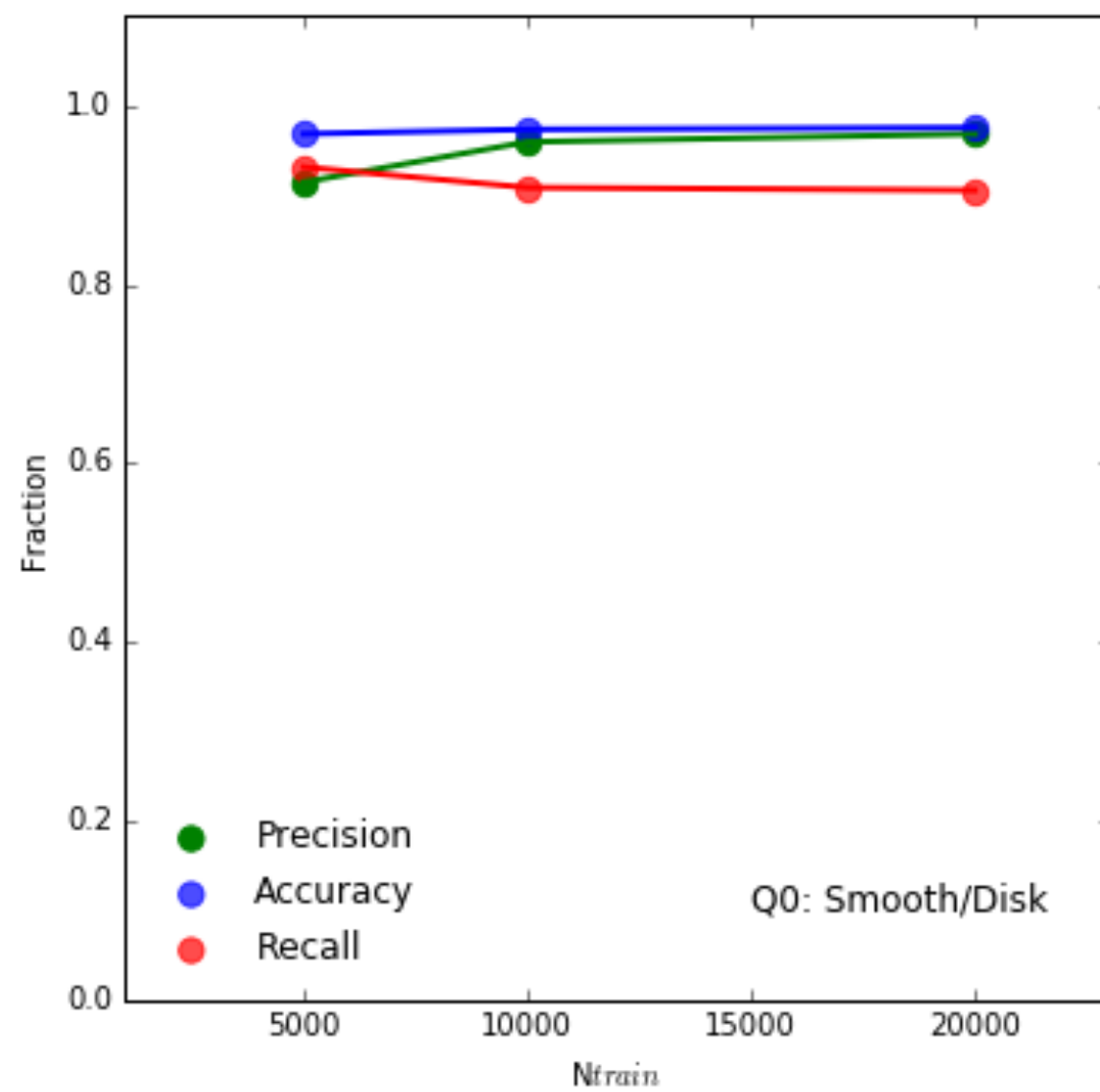


Galaxy ZOO like classifications for for the entire sample

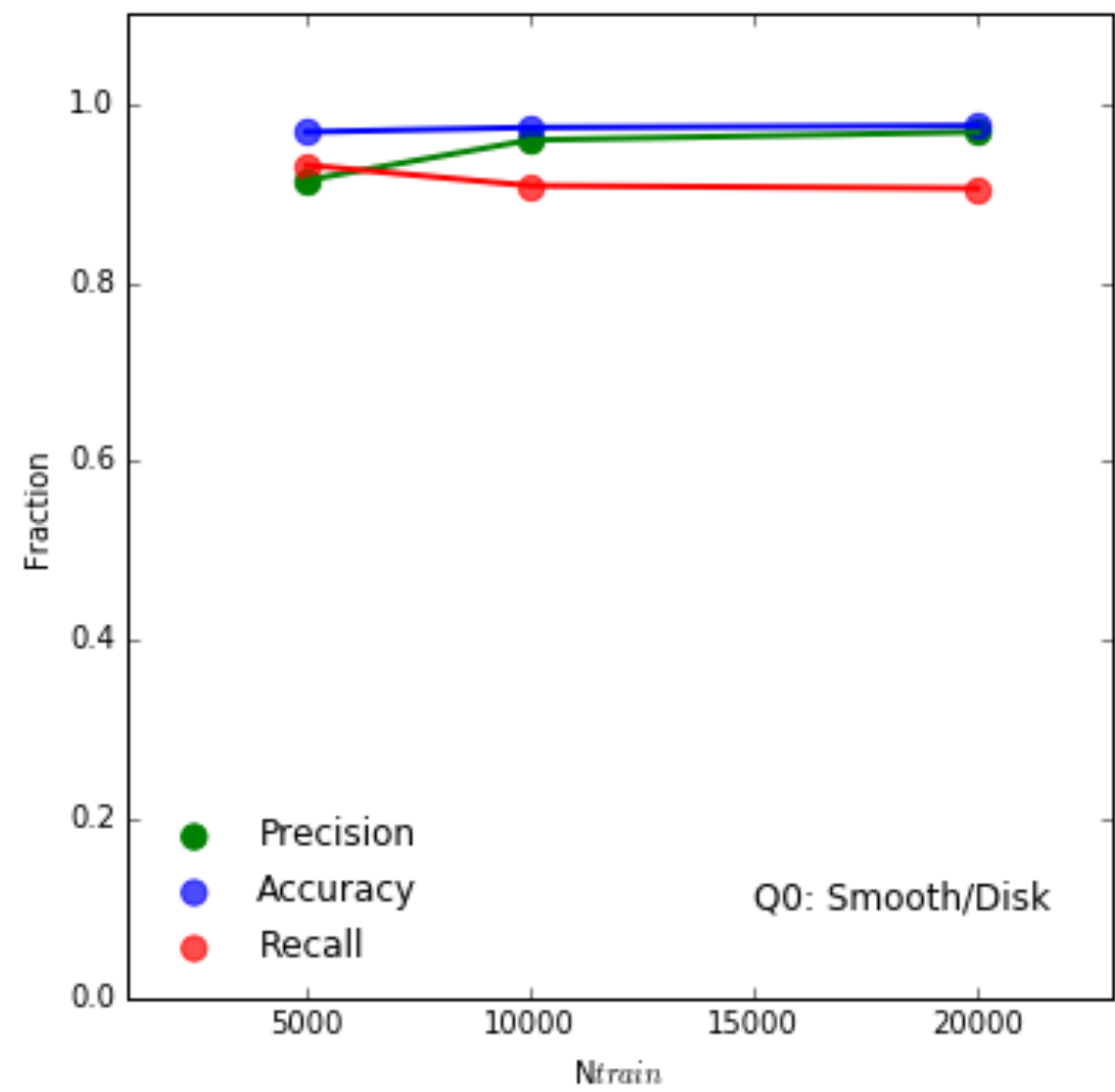
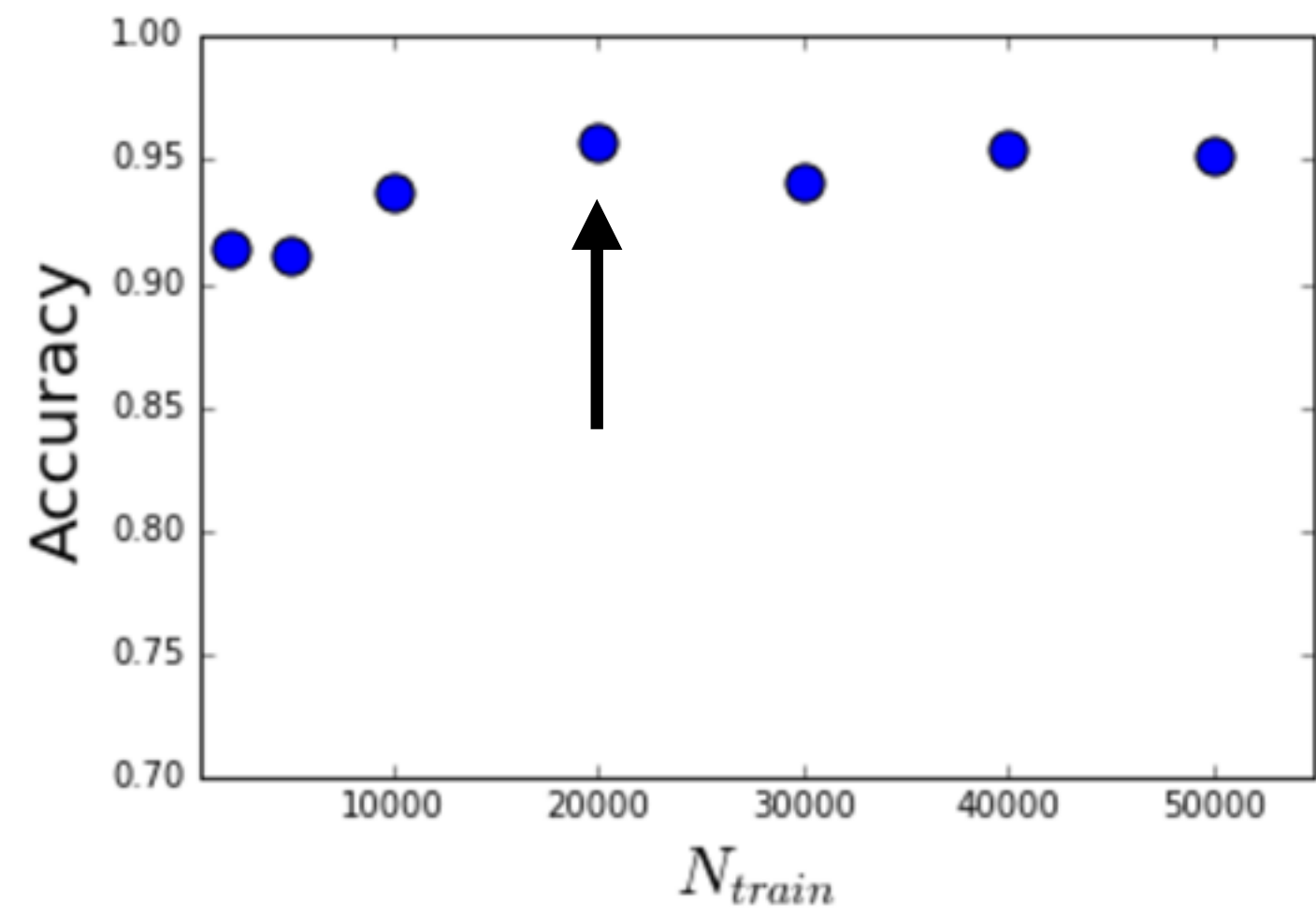
**N??**

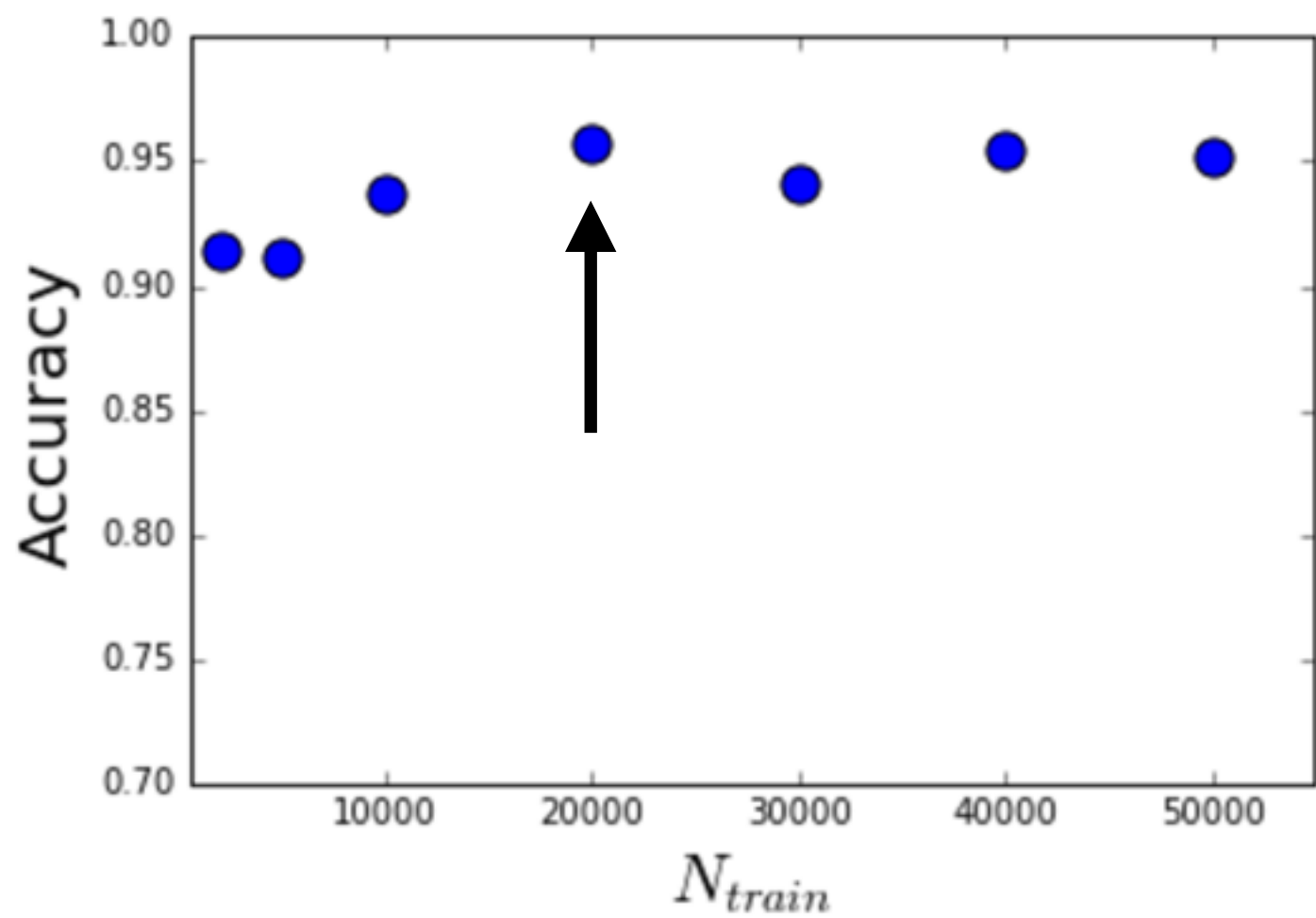
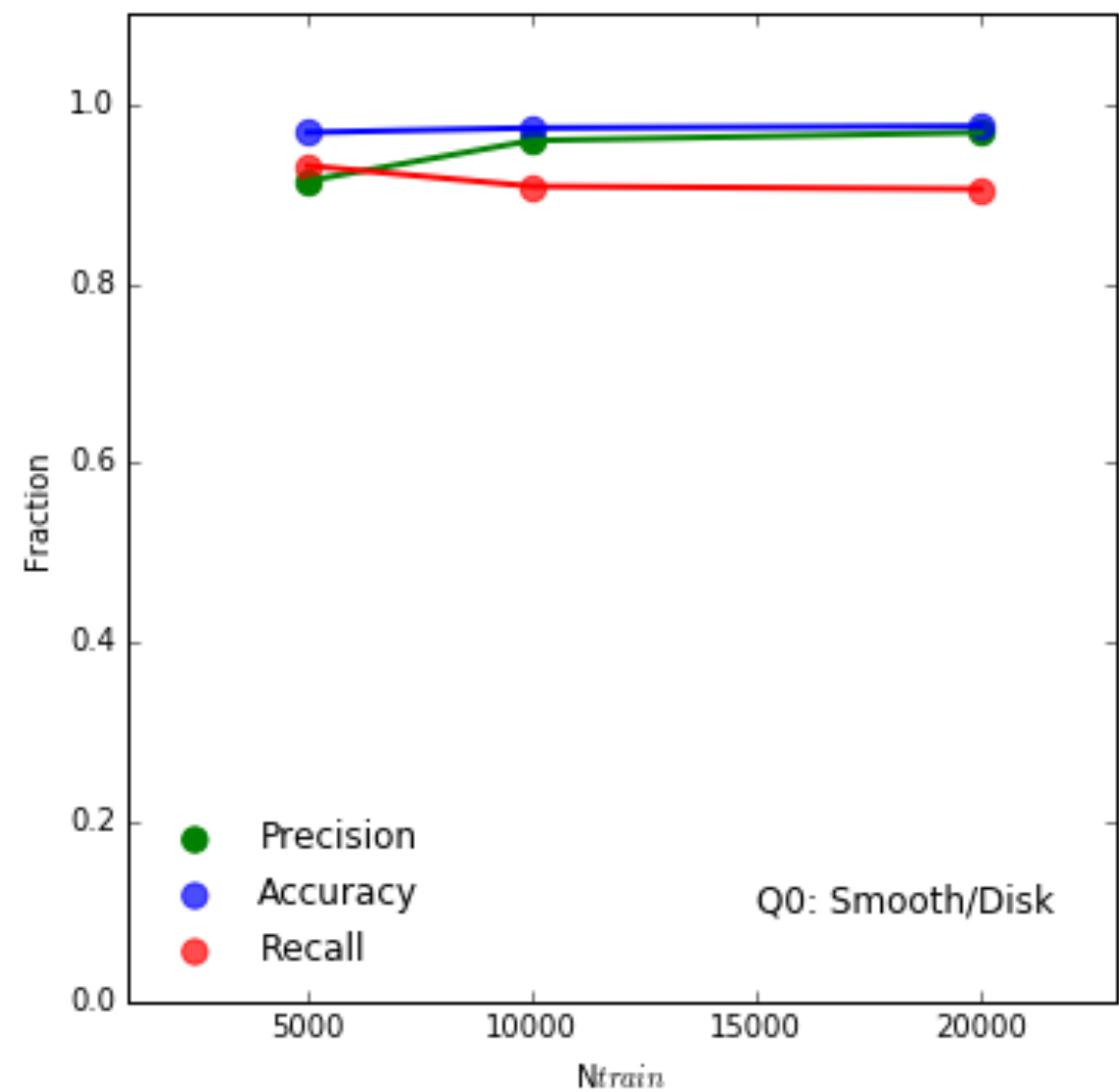
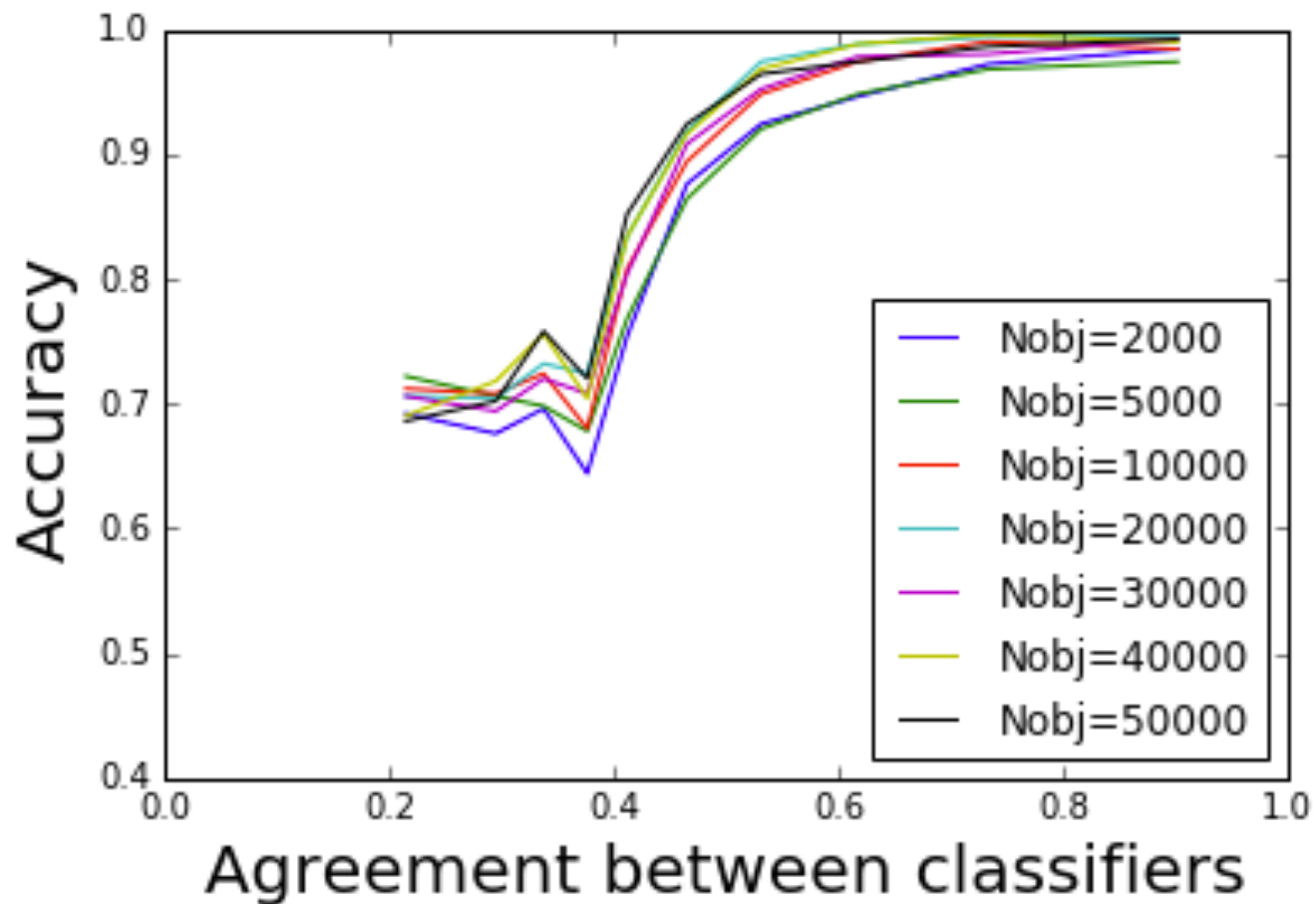
**Knowledge transfer?**

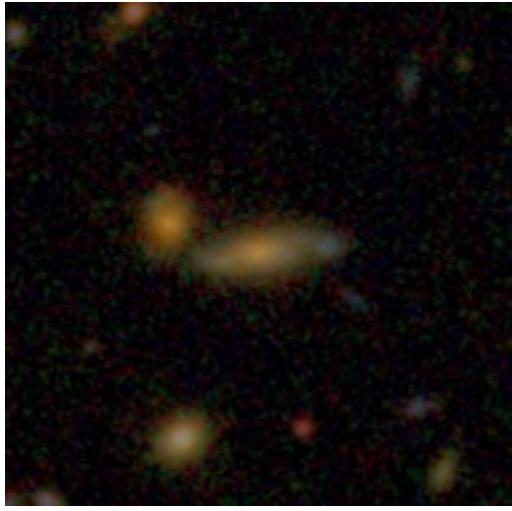
**Human classifications**







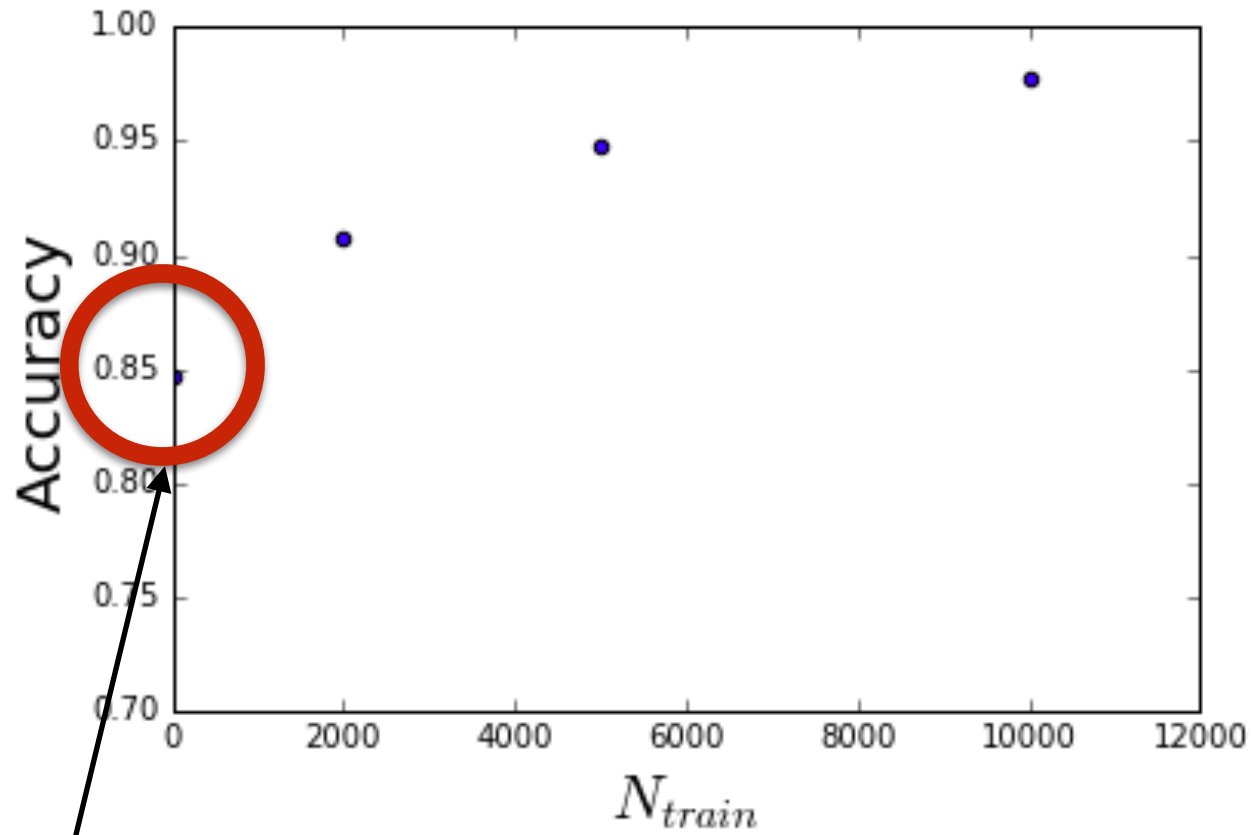




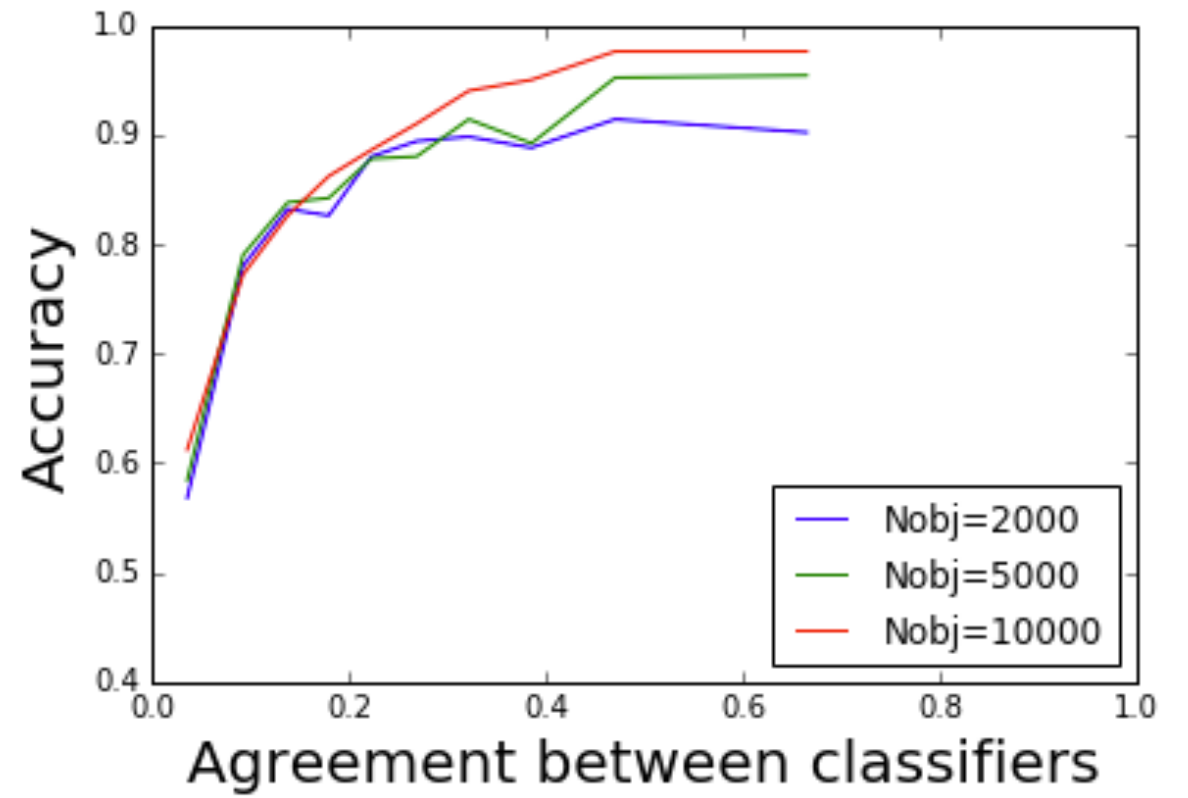
**CANDELS**



DEEP-LEARNING  
MACHINE  
TUNED  
FOR  
SDSS



No training



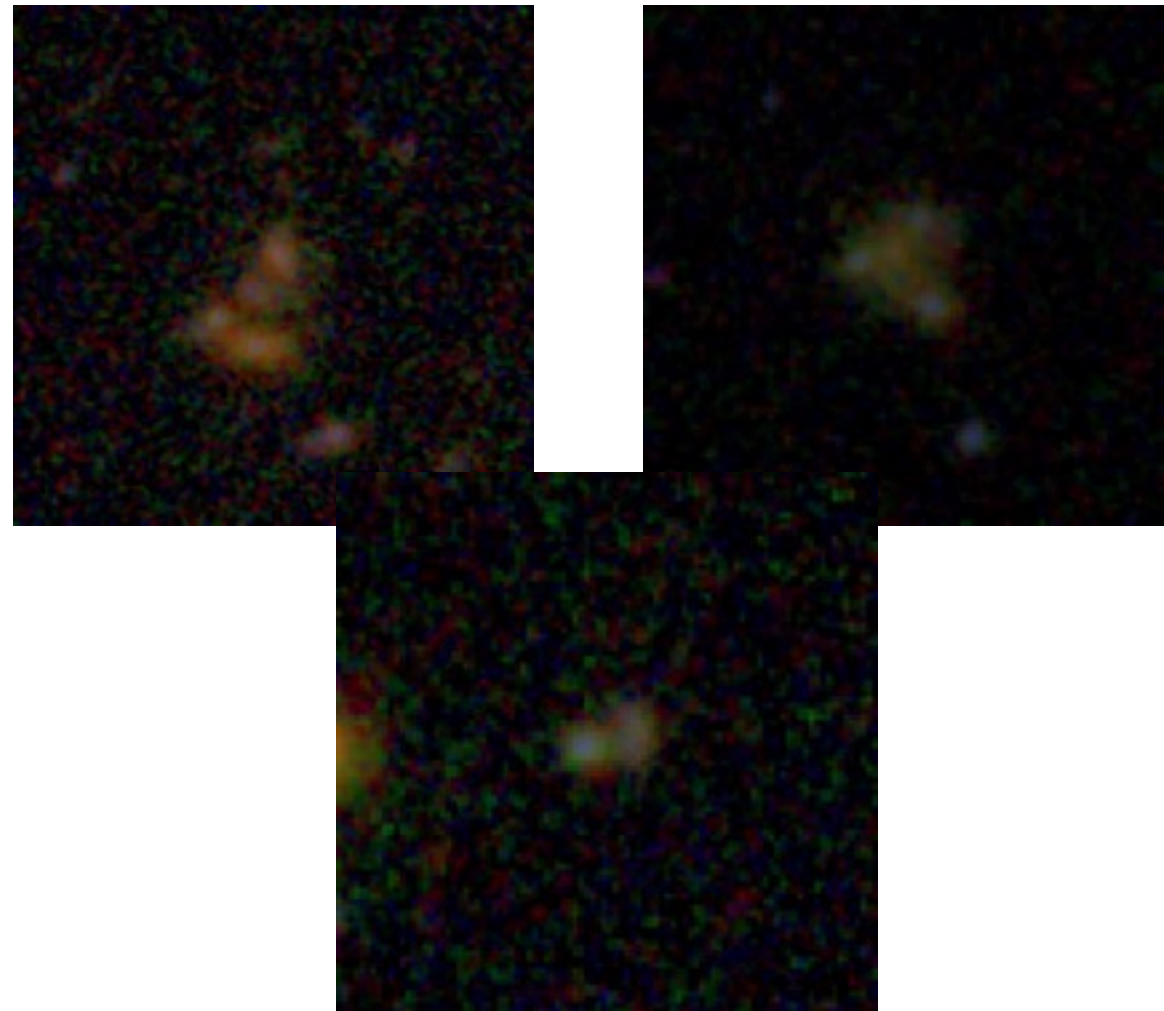
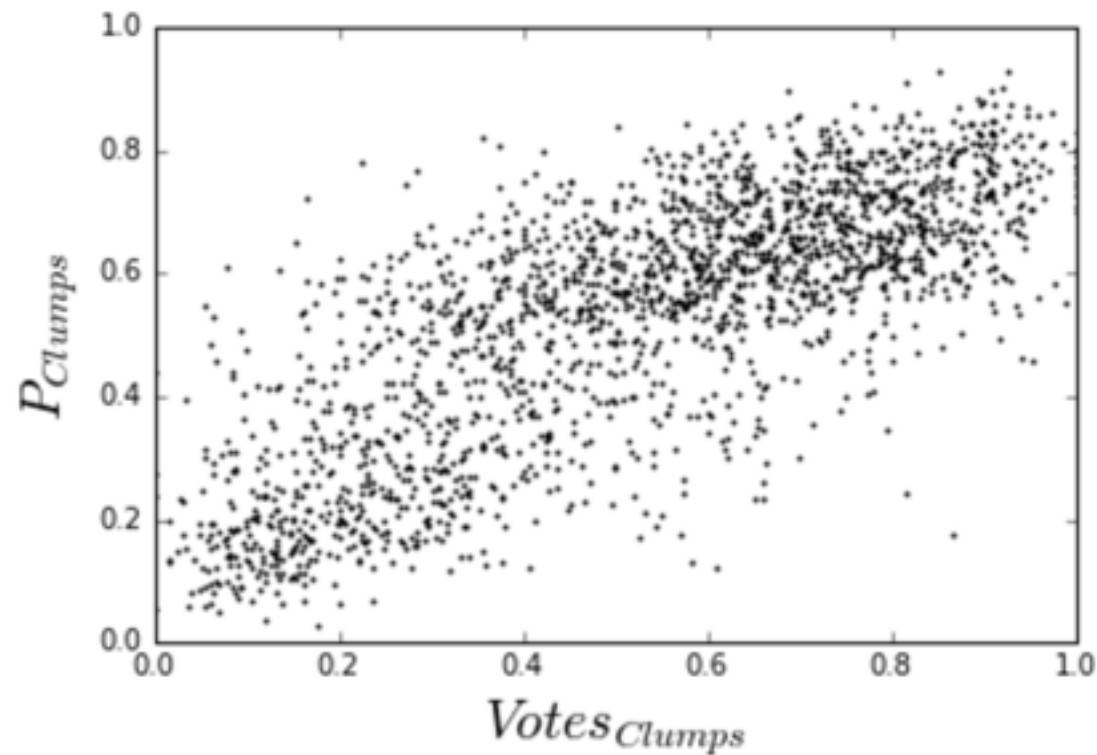
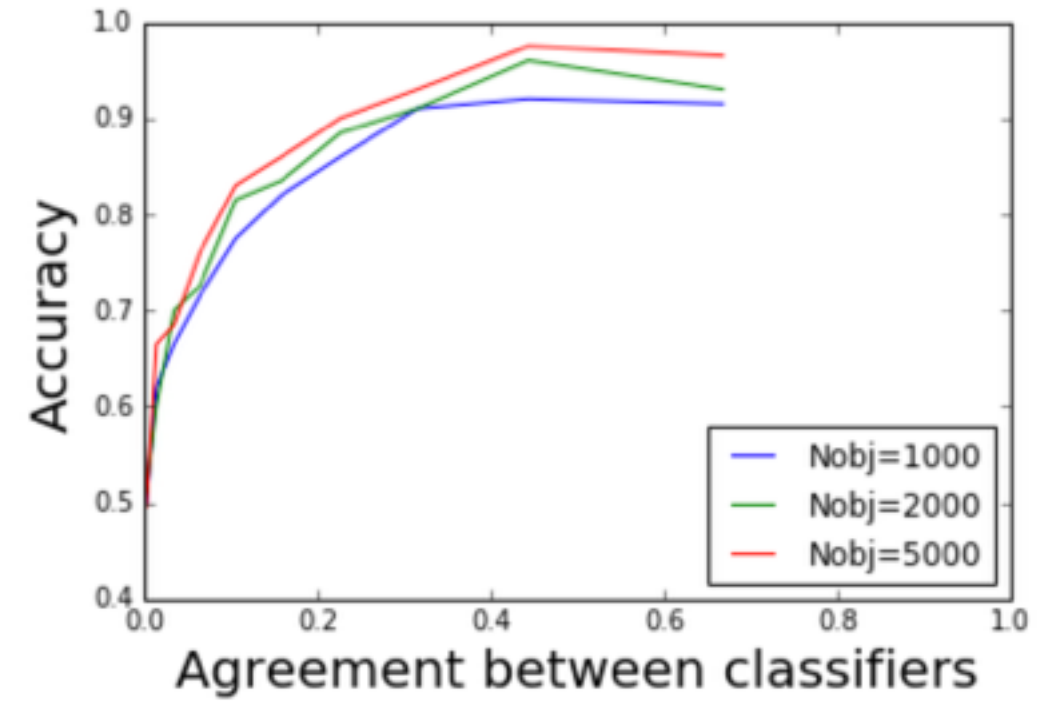
# CLUMPS

DL machine  
trained on  
CANDELS  
first question

update the weights

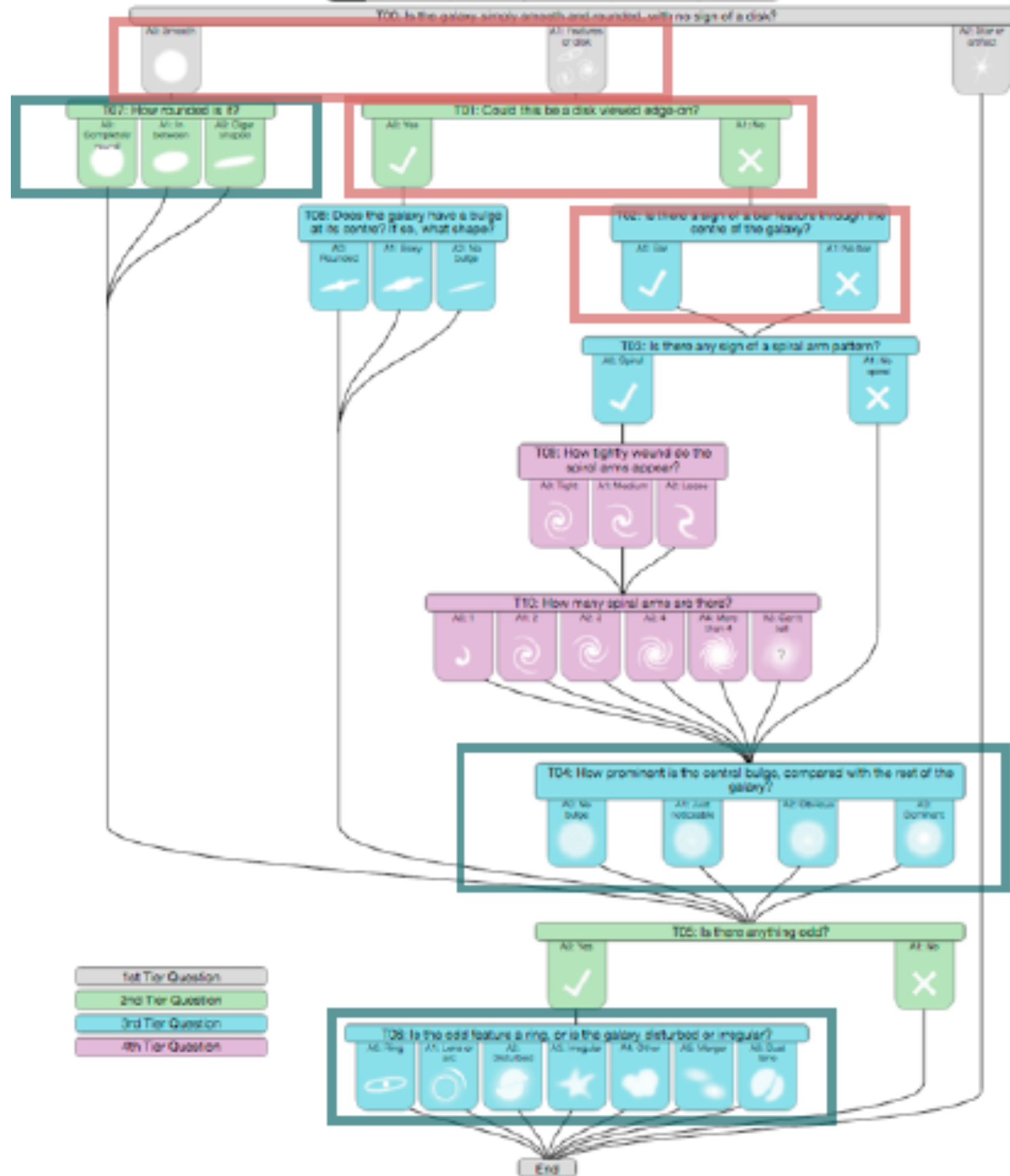


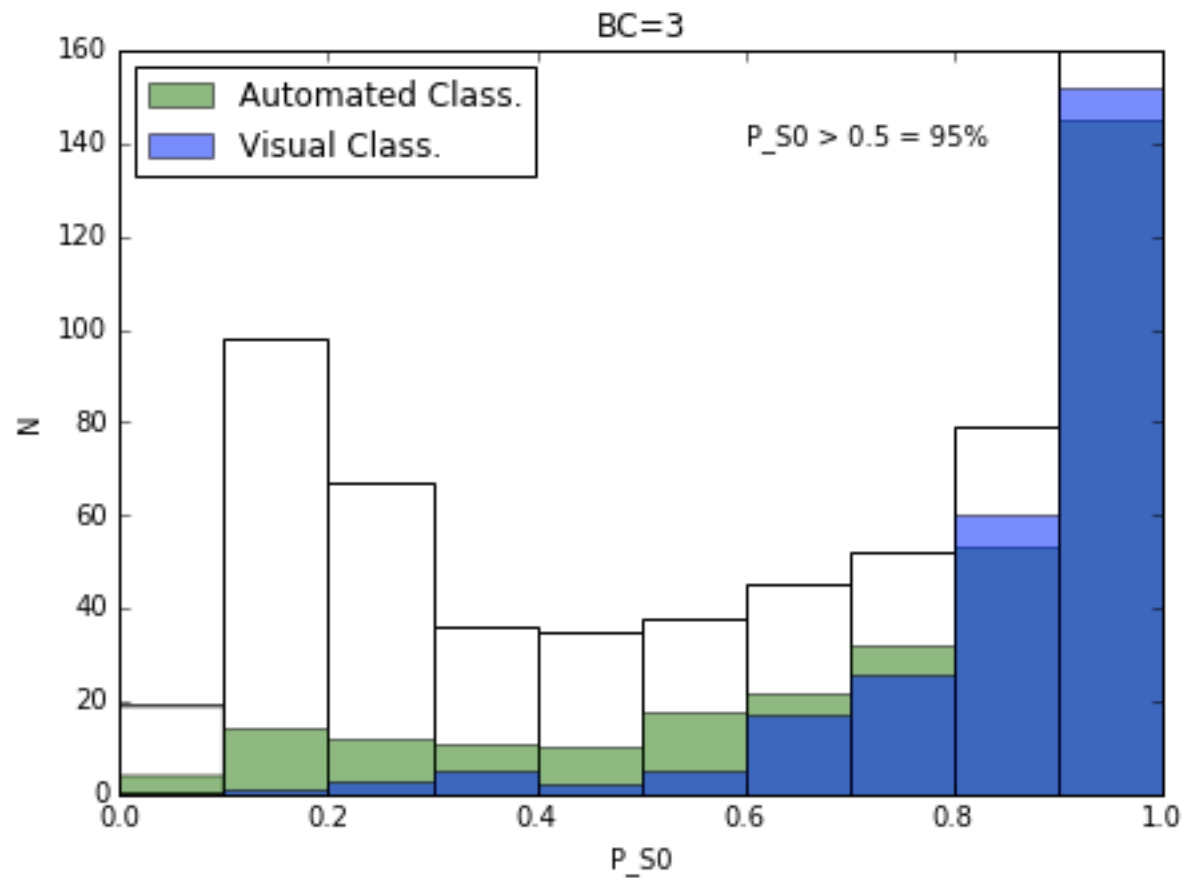
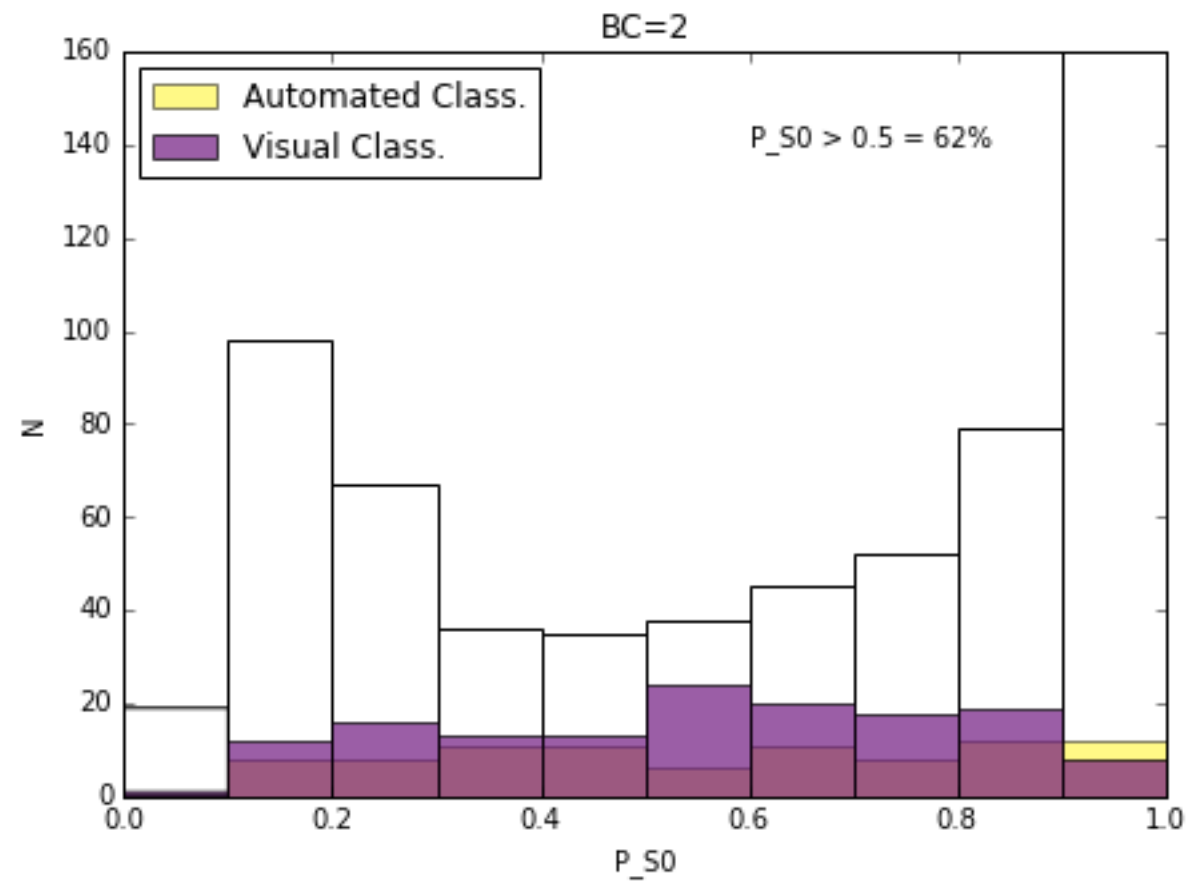
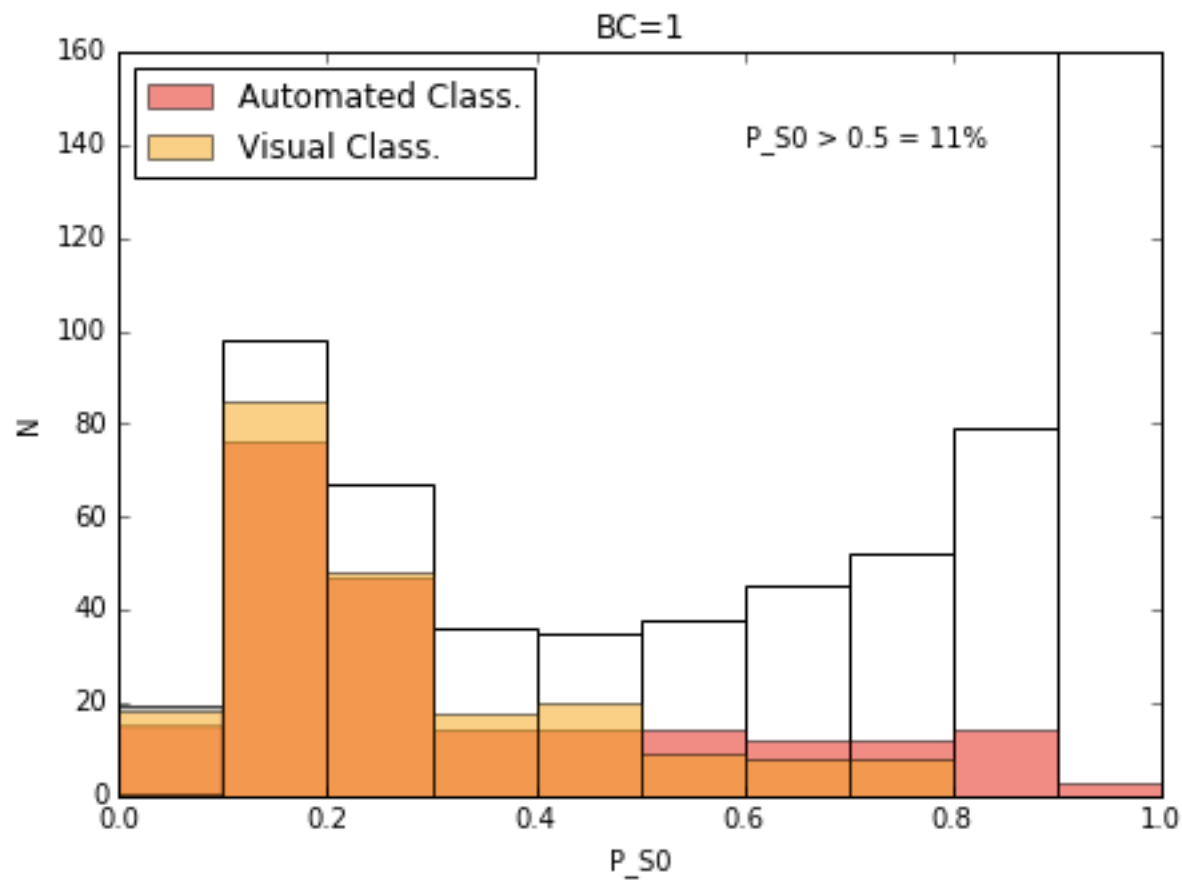
question: clumps  
(2nd level question)



# Galaxy Zoo Decision Trees

QZ 3 QZ 3 Hubble QZ 4 Sloan QZ 4 Galex QZ 4 UKIDSS QZ 4 Fering





Comparison with Cheng,  
Faber+11

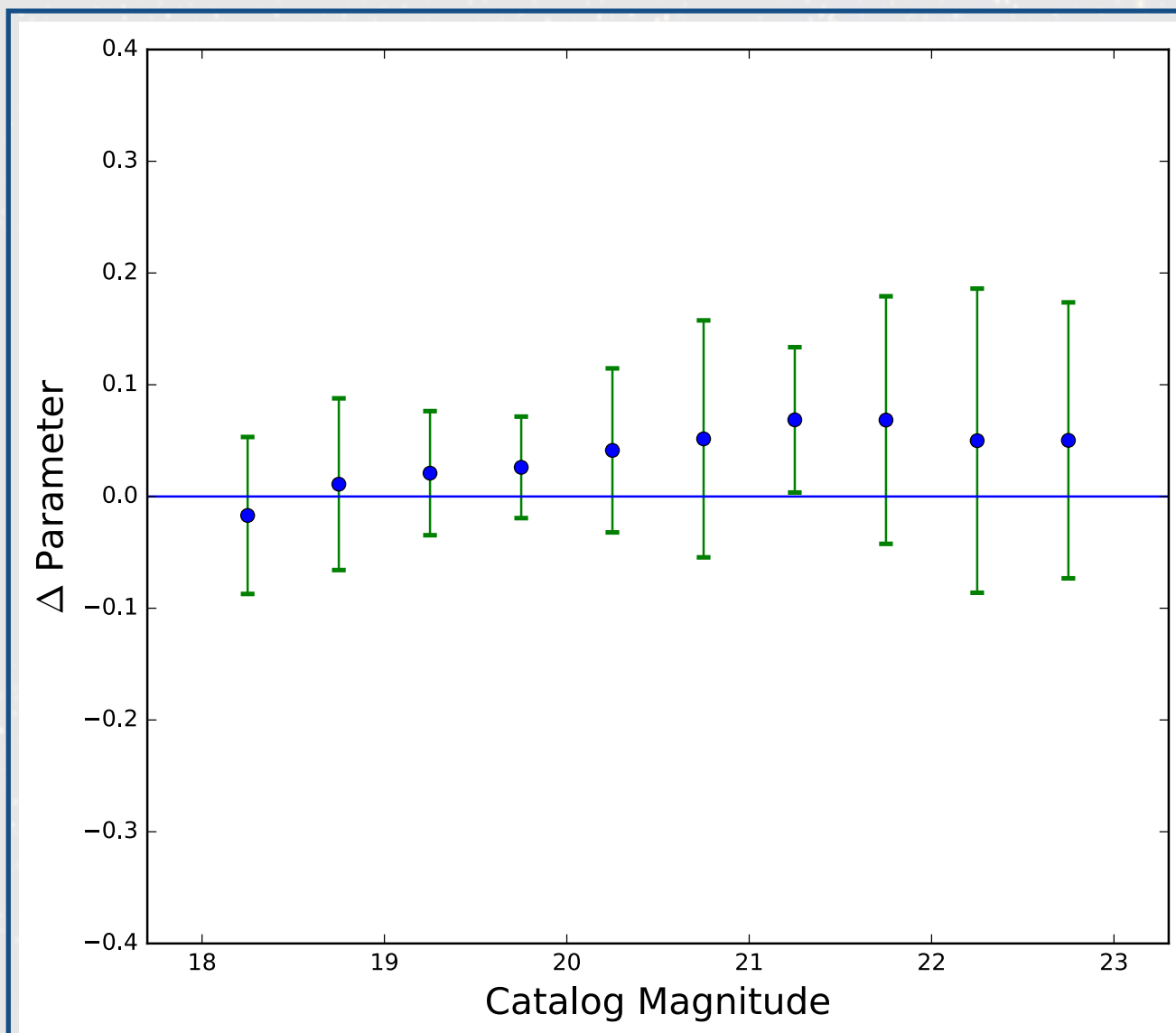


# Group #2: Quantitative measurements

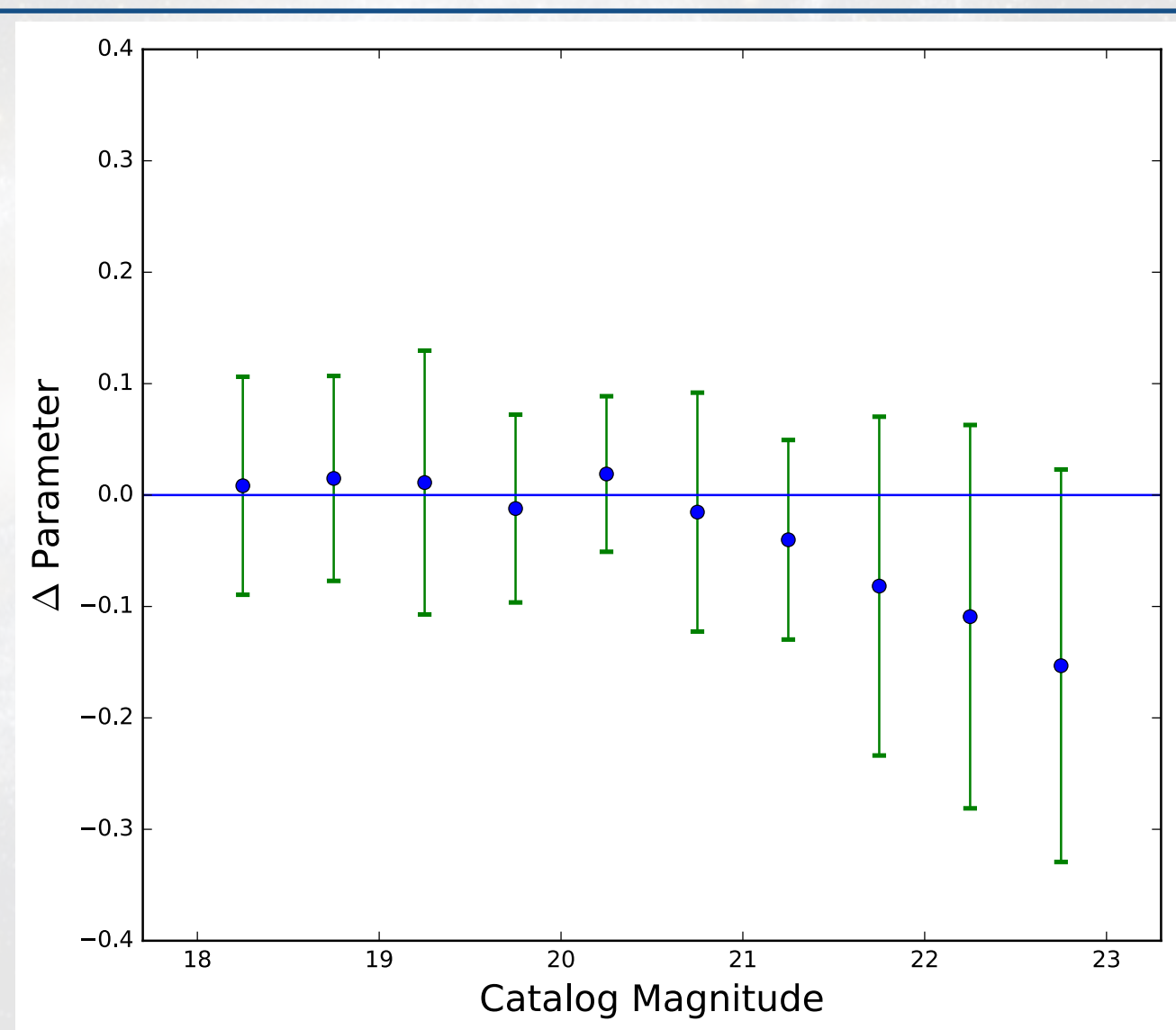
# Predictions on Simulated Data

5000 stamps

Magnitude



DNN

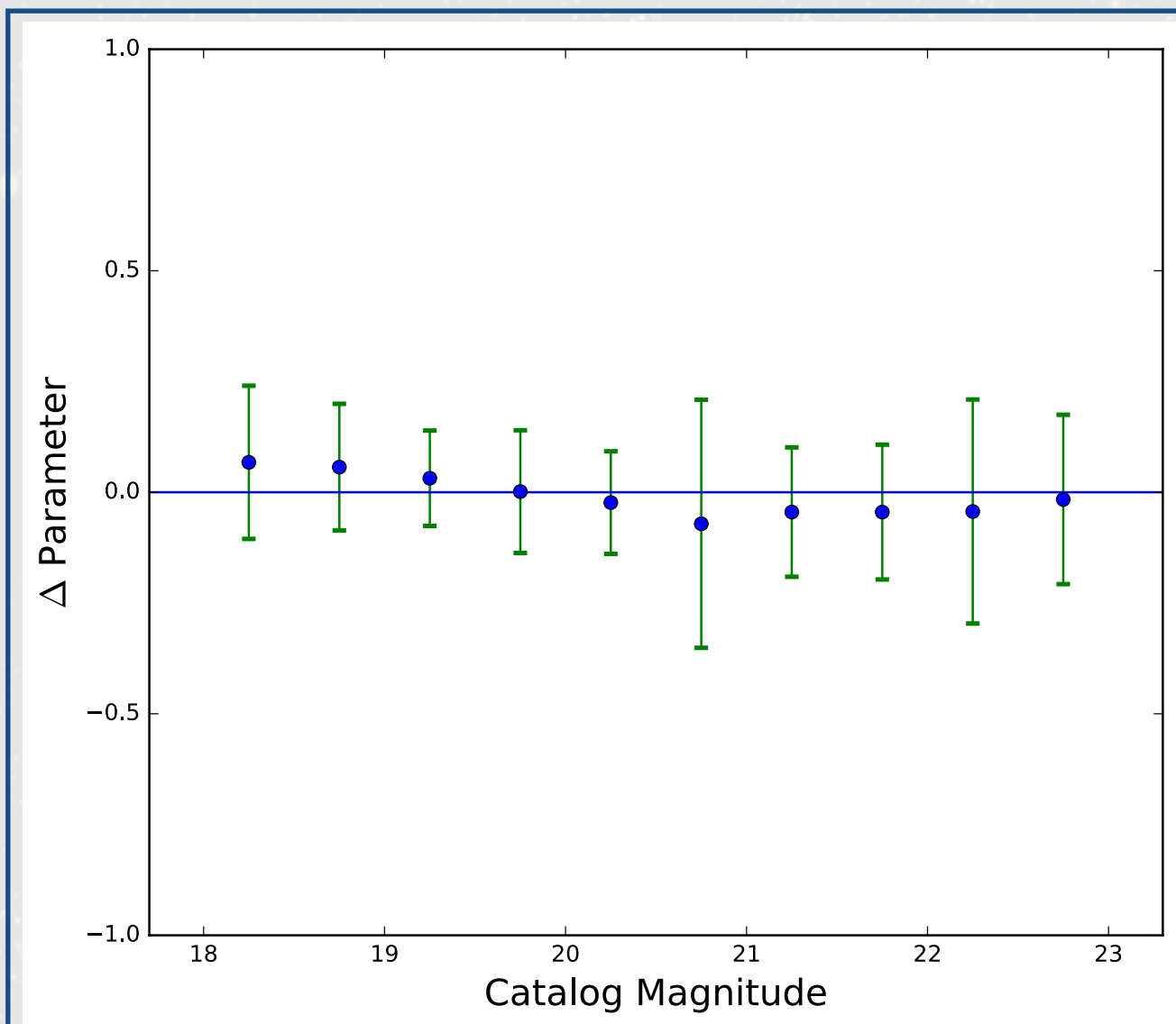


GALFIT-SExtractor

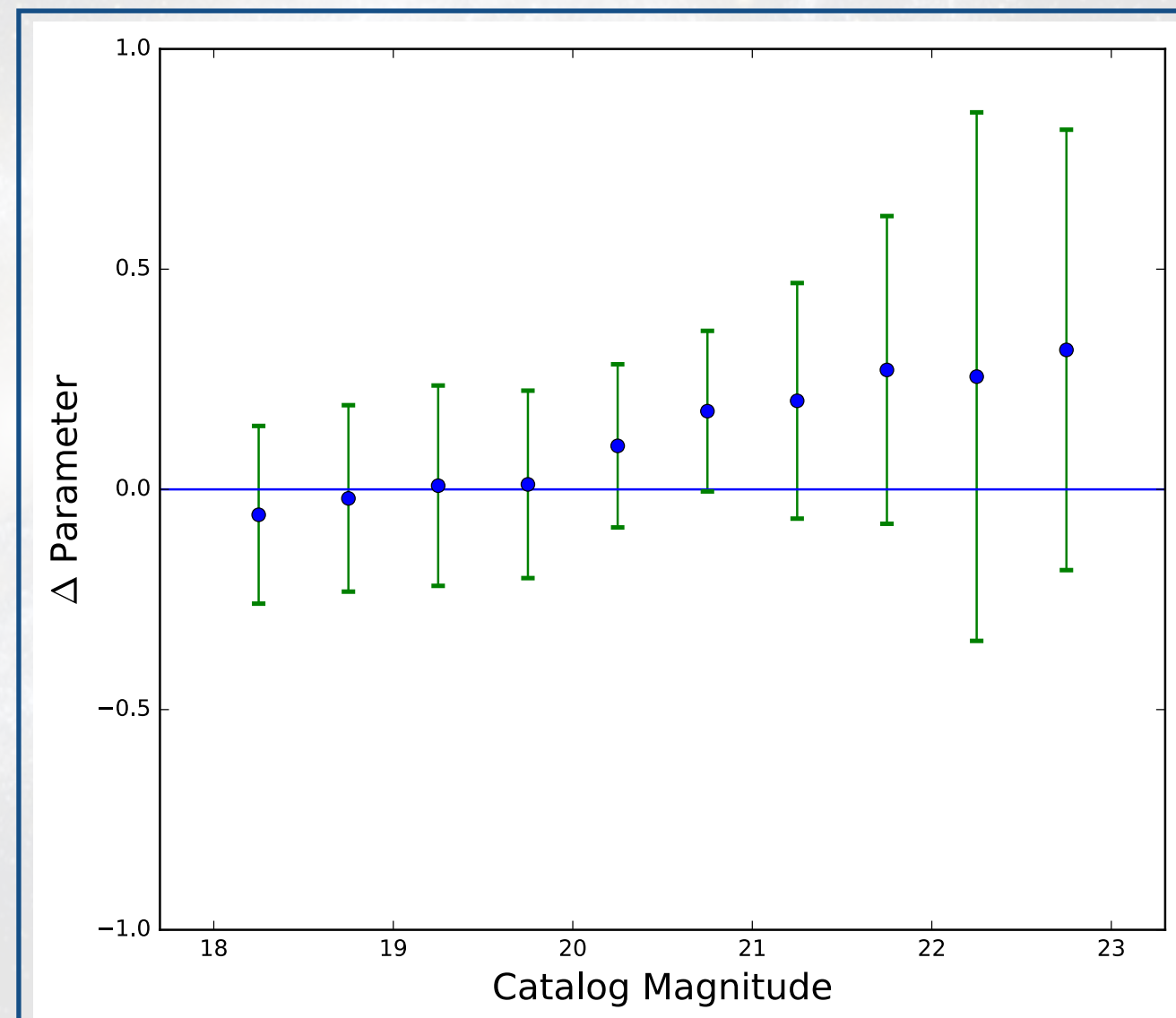
# Predictions on Simulated Data

5000 stamps

Sersic index



DNN



GALFIT-SExtractor



# Summary of predictions on simulation

$R^2$ simulated data			
Parameter	Architecture 1	Architecture 2	GALFIT
Magnitude	0.947	0.995	0.986
Radius	0.892	0.955	0.738
Sérsic index	0.887	0.348	0.292
Ellipticity	0.755	0.603	0.896
Position Angle	0.941	nc	0.825

coefficient of determination

$$R^2 = 1 - \frac{\sum_i (y_i - f_i)^2}{\sum_i (y_i - \bar{y})^2}$$

~5000  
sources fitted



GALFIT/GALAPAGOS

Time: ~ 4 Hours



CNN (ONCE TRAINED)

Time: < 3 seconds



# Summary of Predictions on Real Data (after domain adaptation)

Parameter	$R^2$ Real data		
	Before TL	After TL	2 GALFIT
Magnitude	0.788	0.982	0.985
Radius	-1.639	0.856	0.860
Sérsic index	-0.768	0.718	0.735
Ellipticity	0.256	0.897	0.904
Position Angle	0.132	0.893	0.863

coefficient of determination

$$R^2 = 1 - \frac{\sum_i (y_i - f_i)^2}{\sum_i (y_i - \bar{y})^2}$$

~3000  
sources fitted

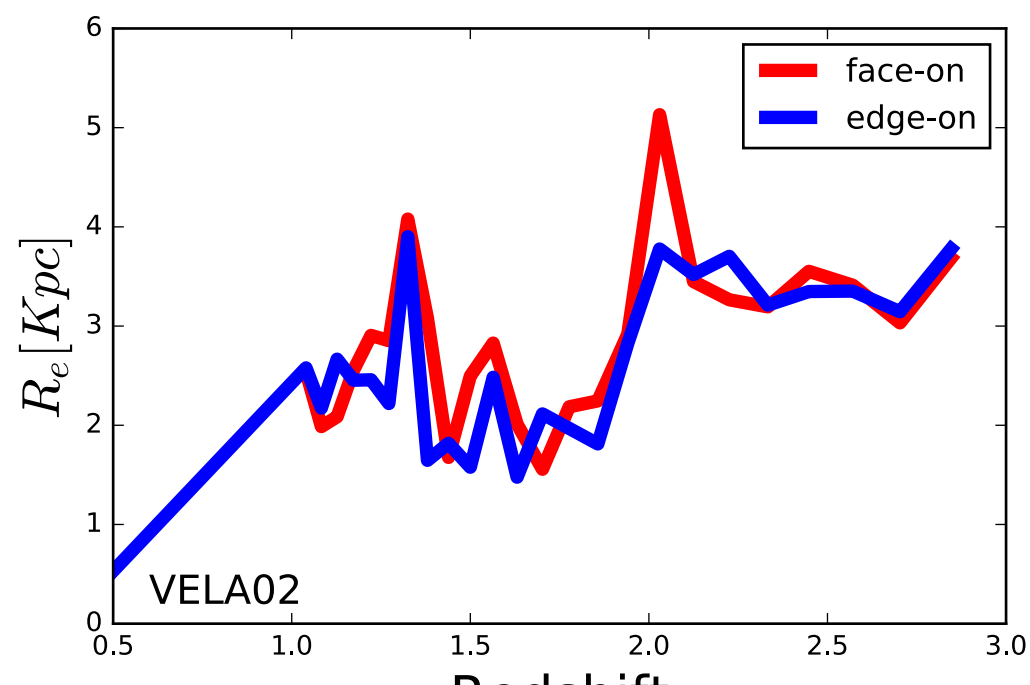
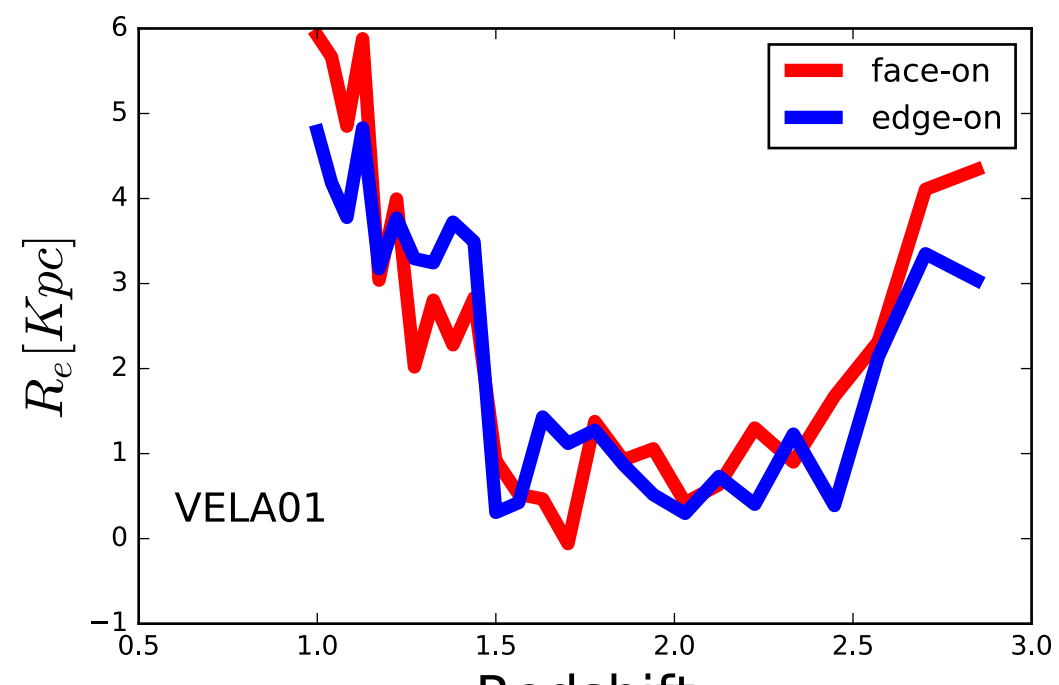
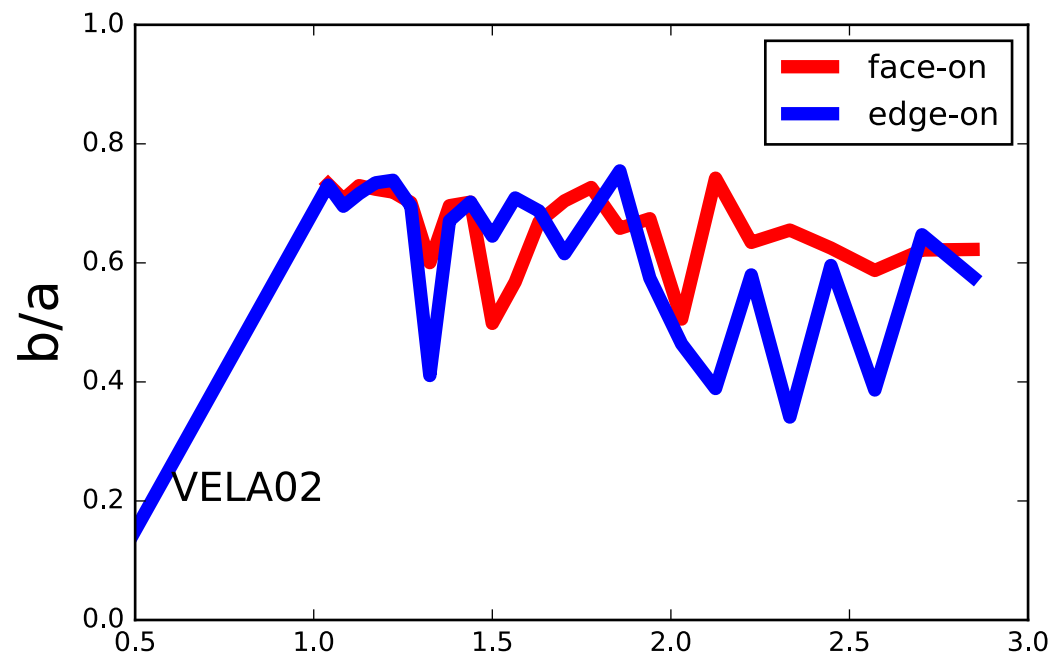
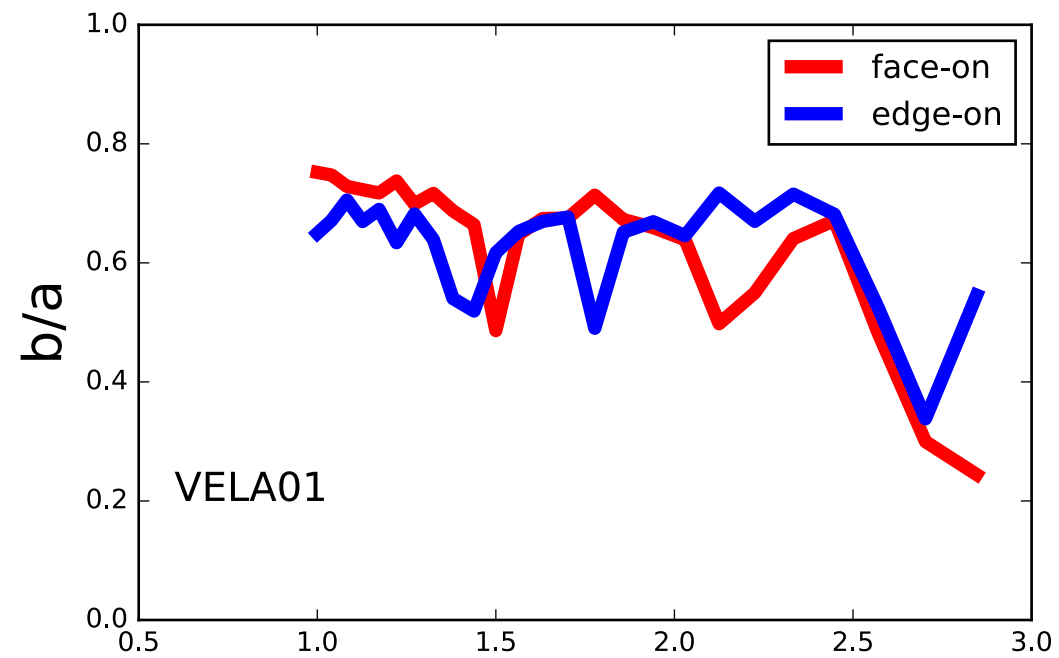
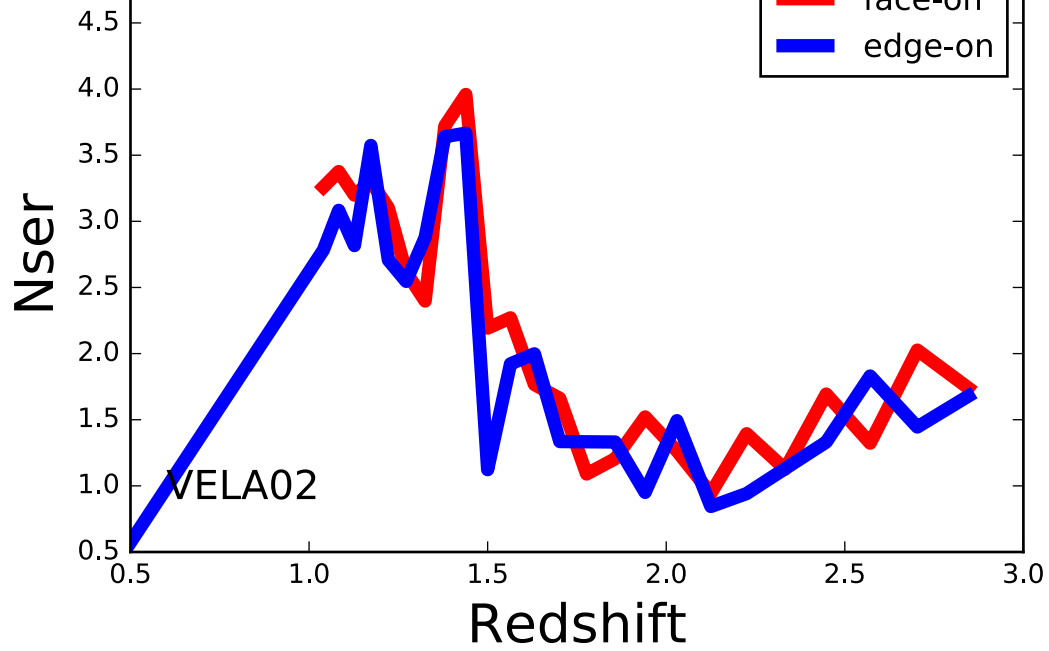
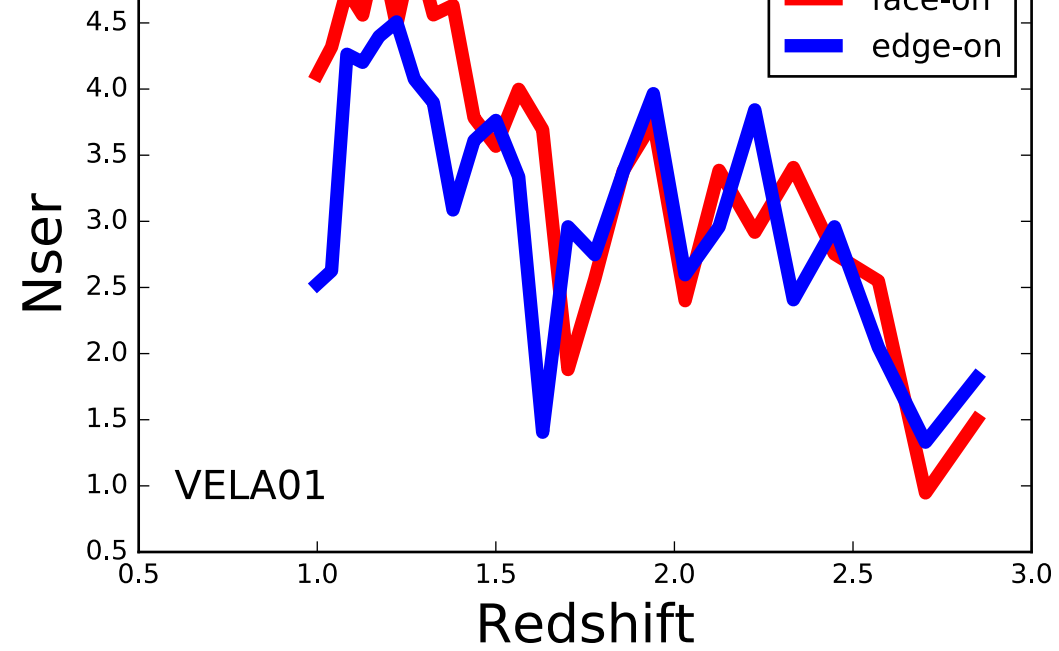
GALFIT/GALAPAGOS

Time: ~ 2h30

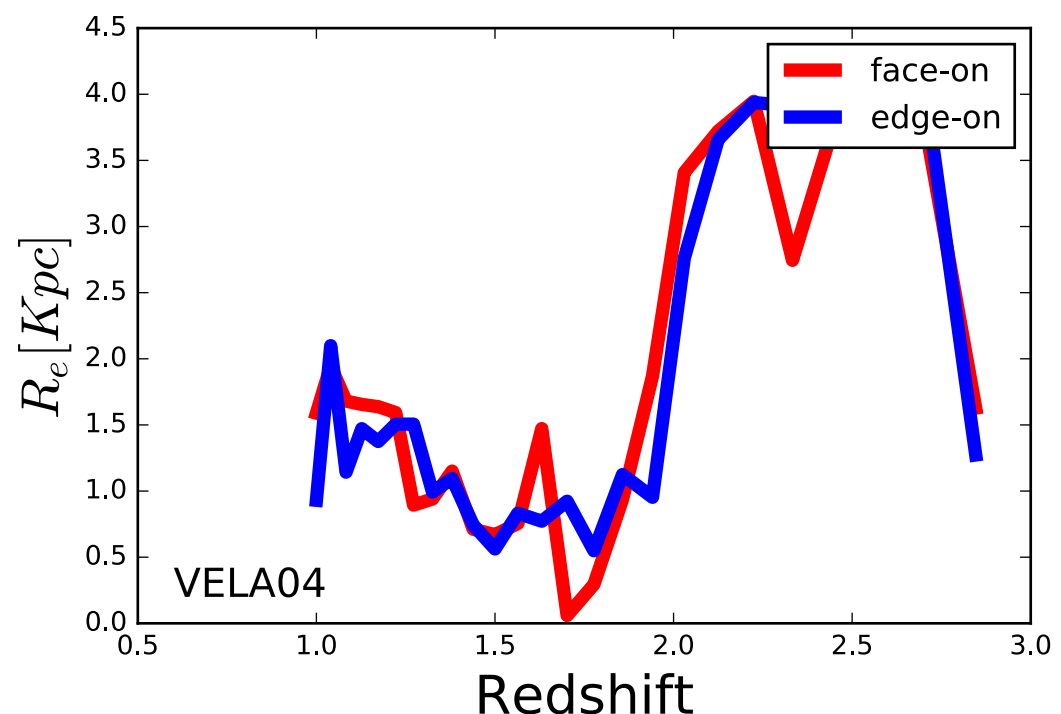
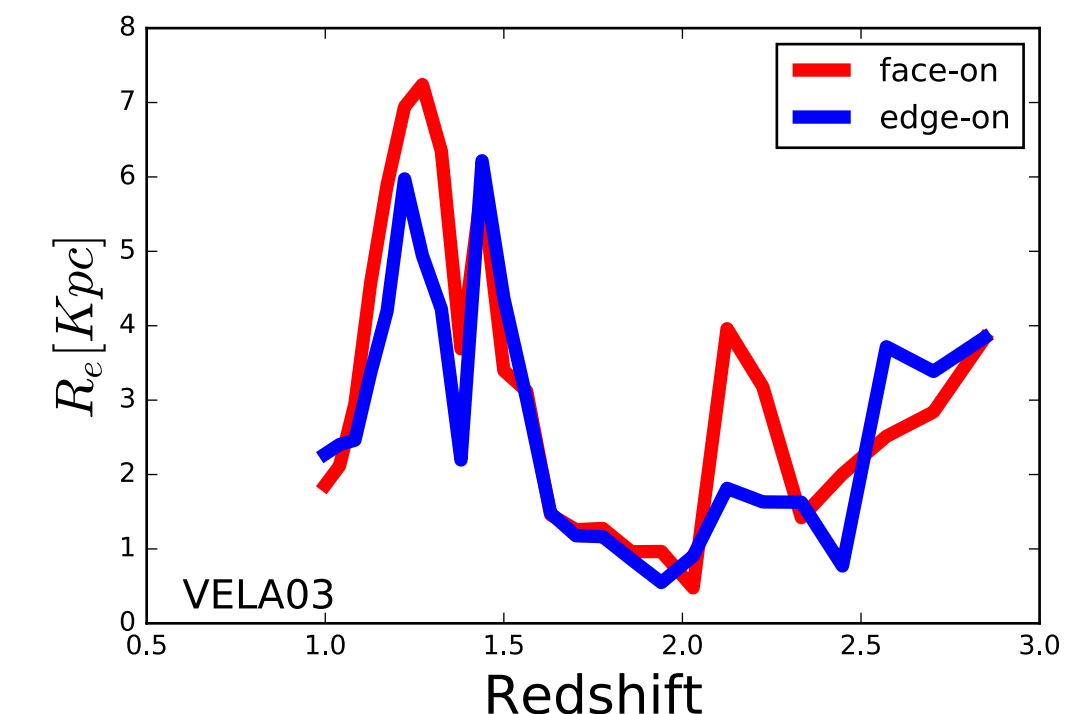
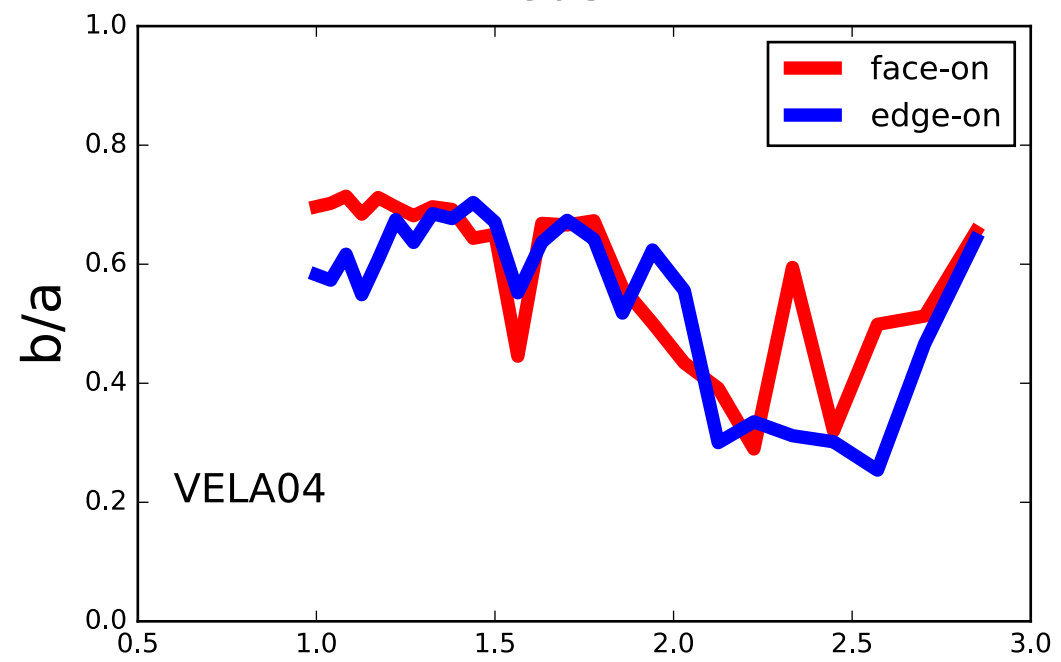
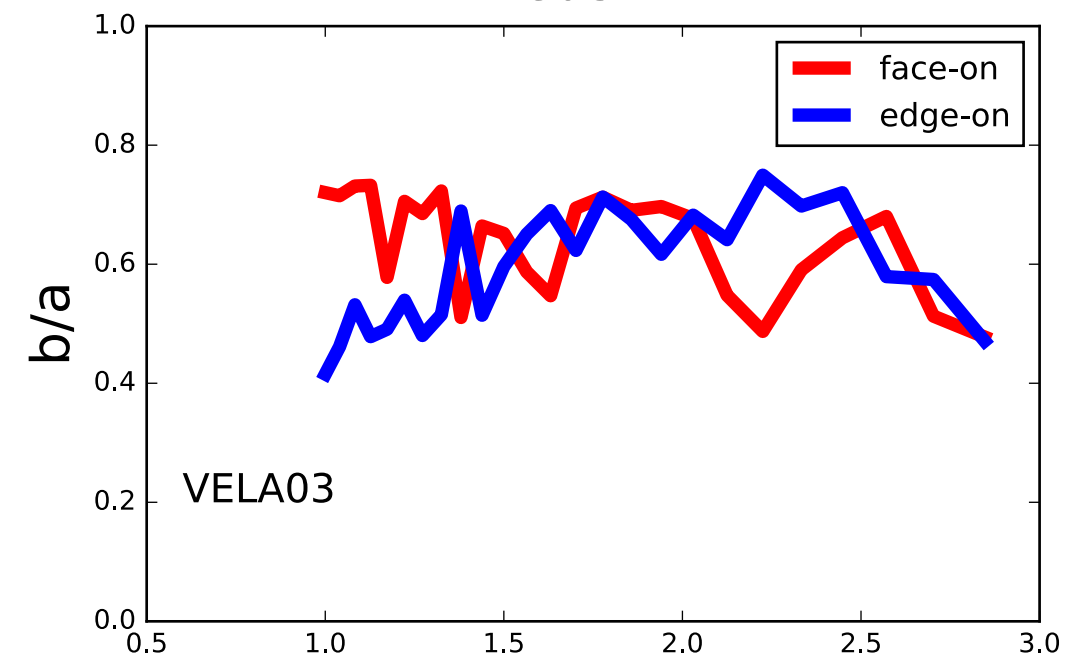
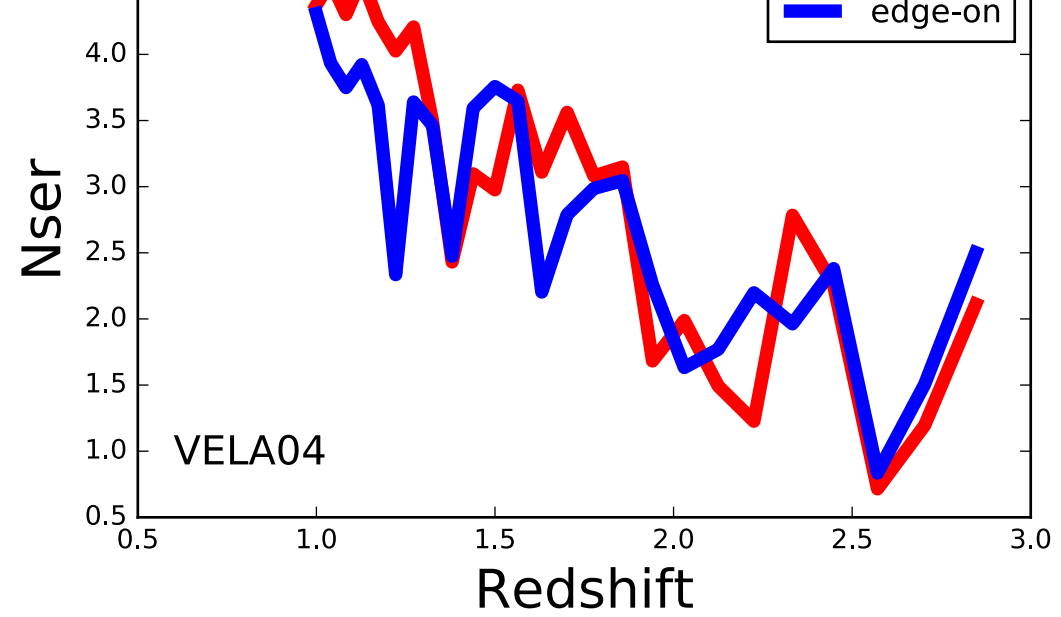
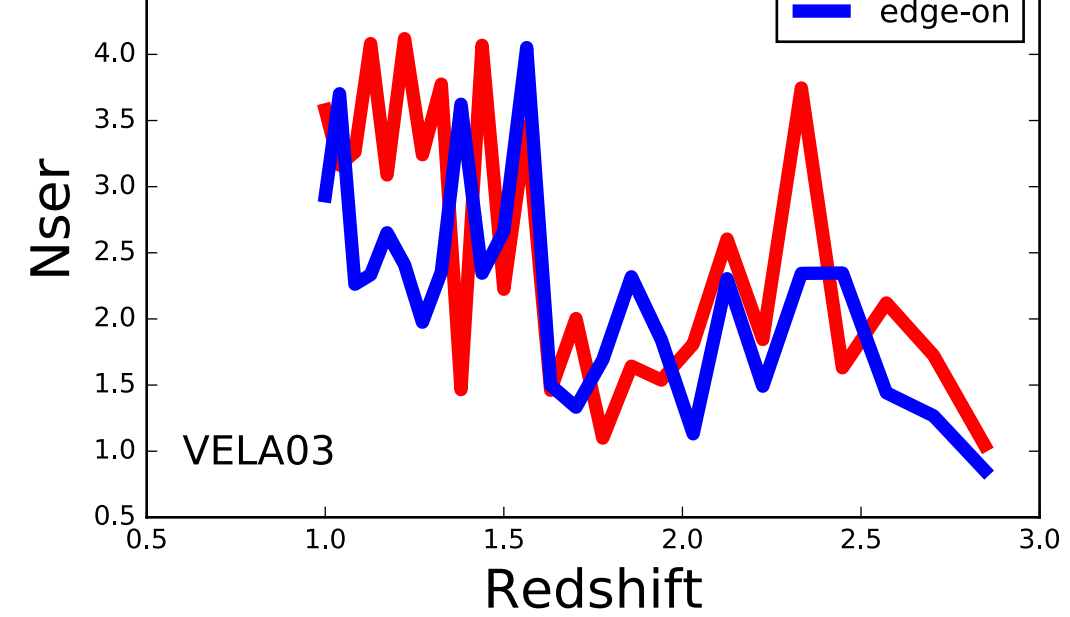
CNN

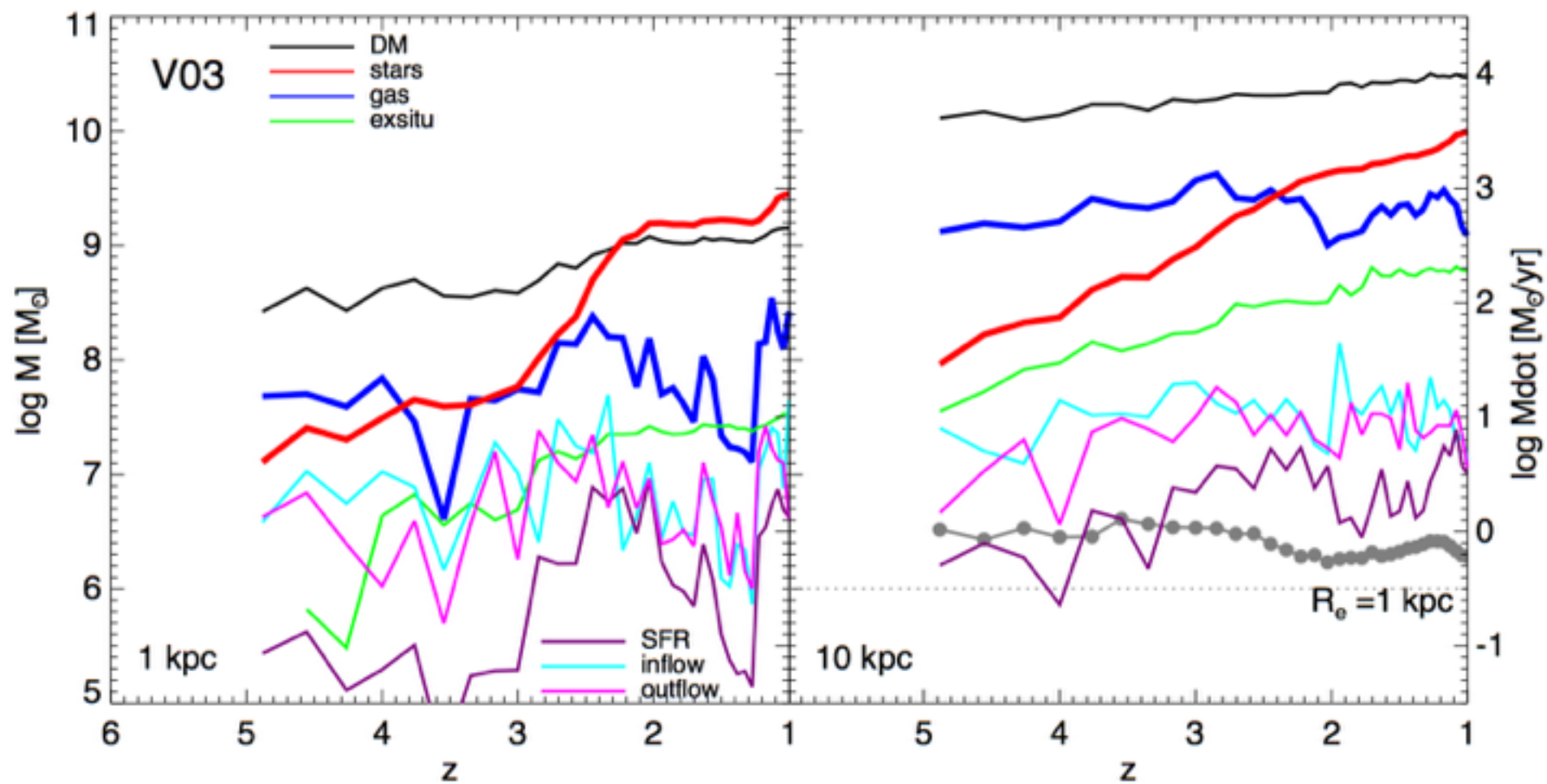
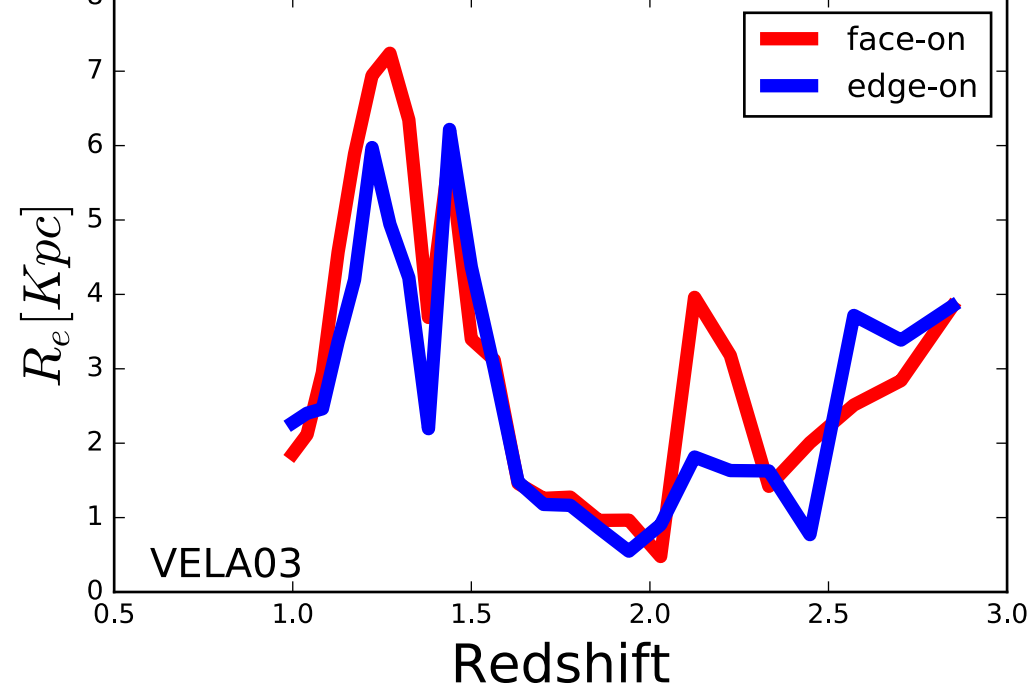
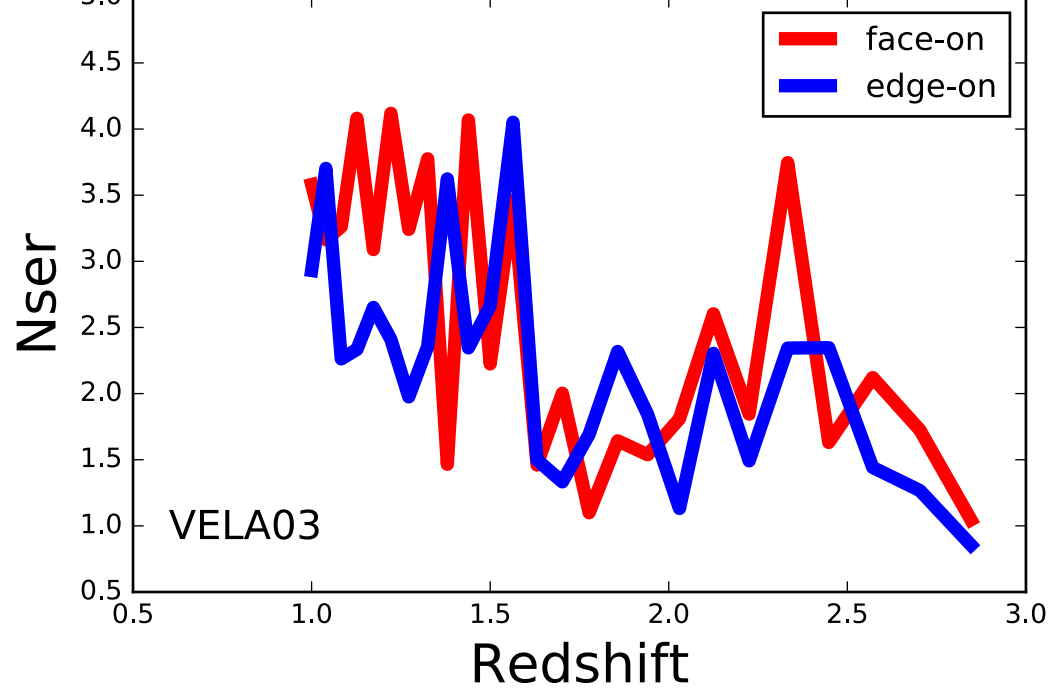
Time: 10 min domain adaptation + 2s

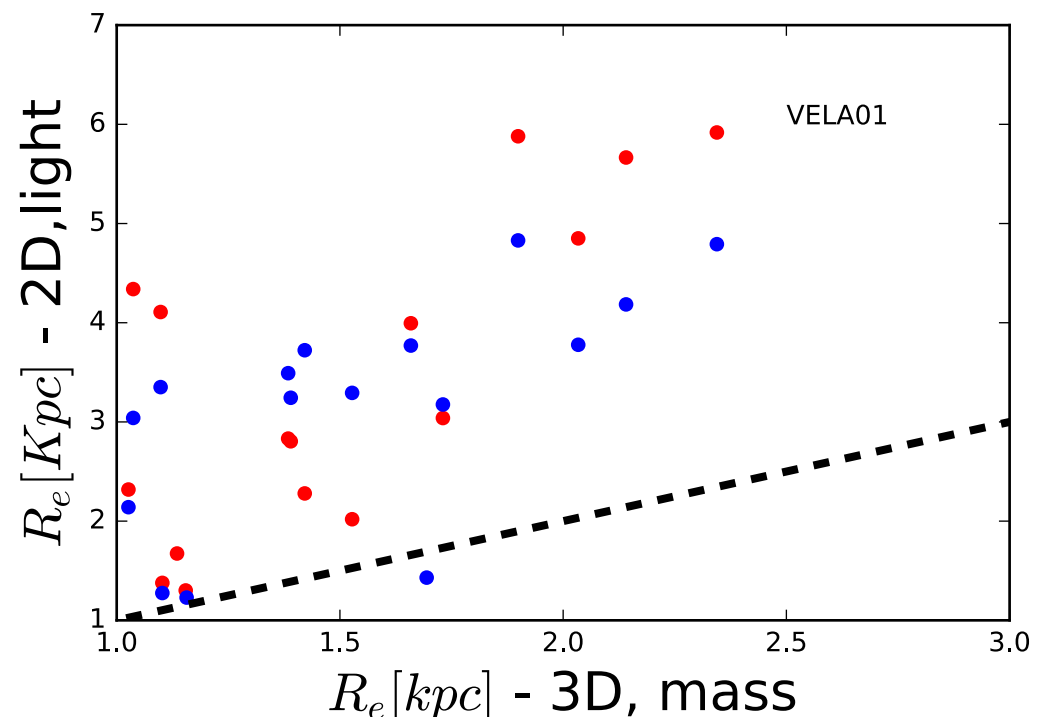
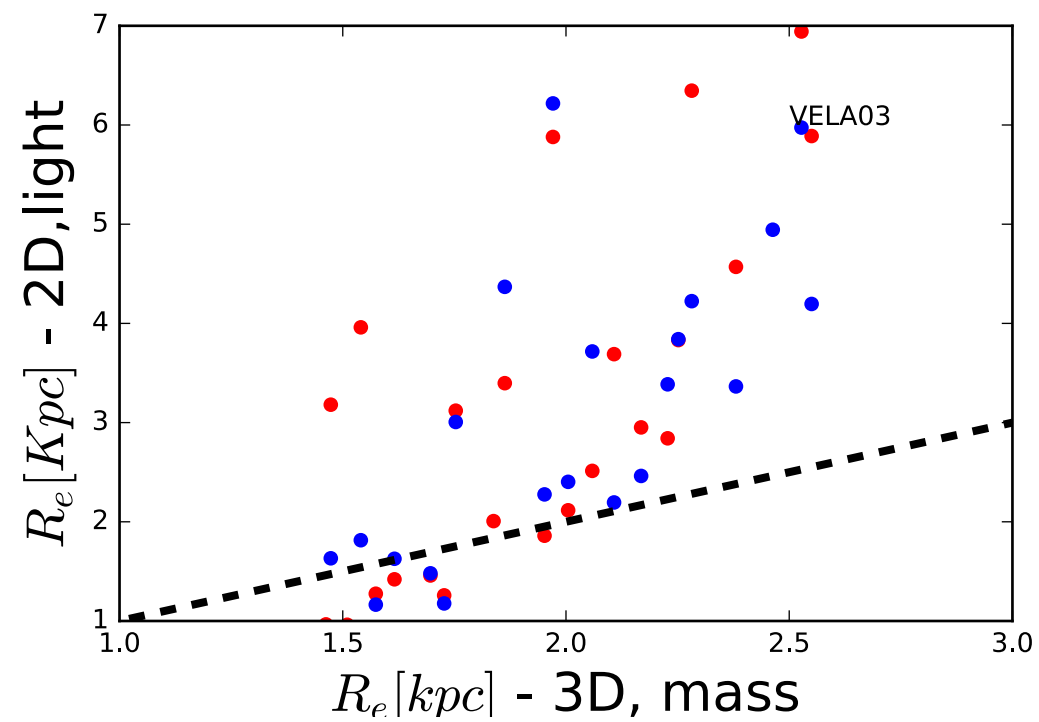
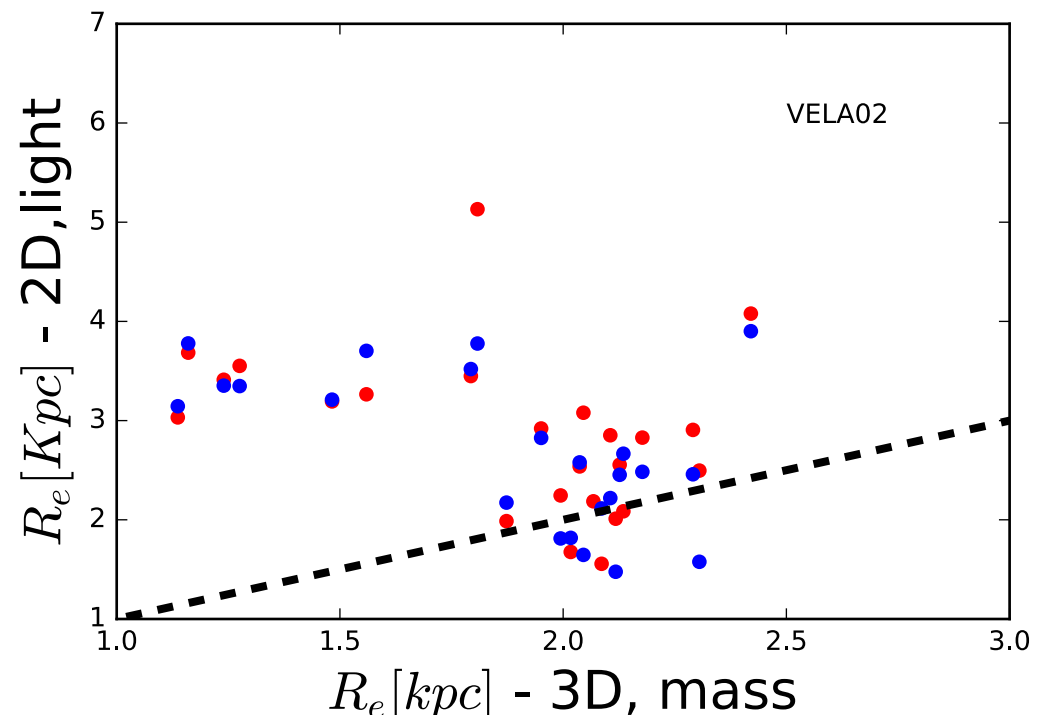












Group #3: Hidden  
observables / correlations