

# **SDSS DR12: Object classification & analysis**

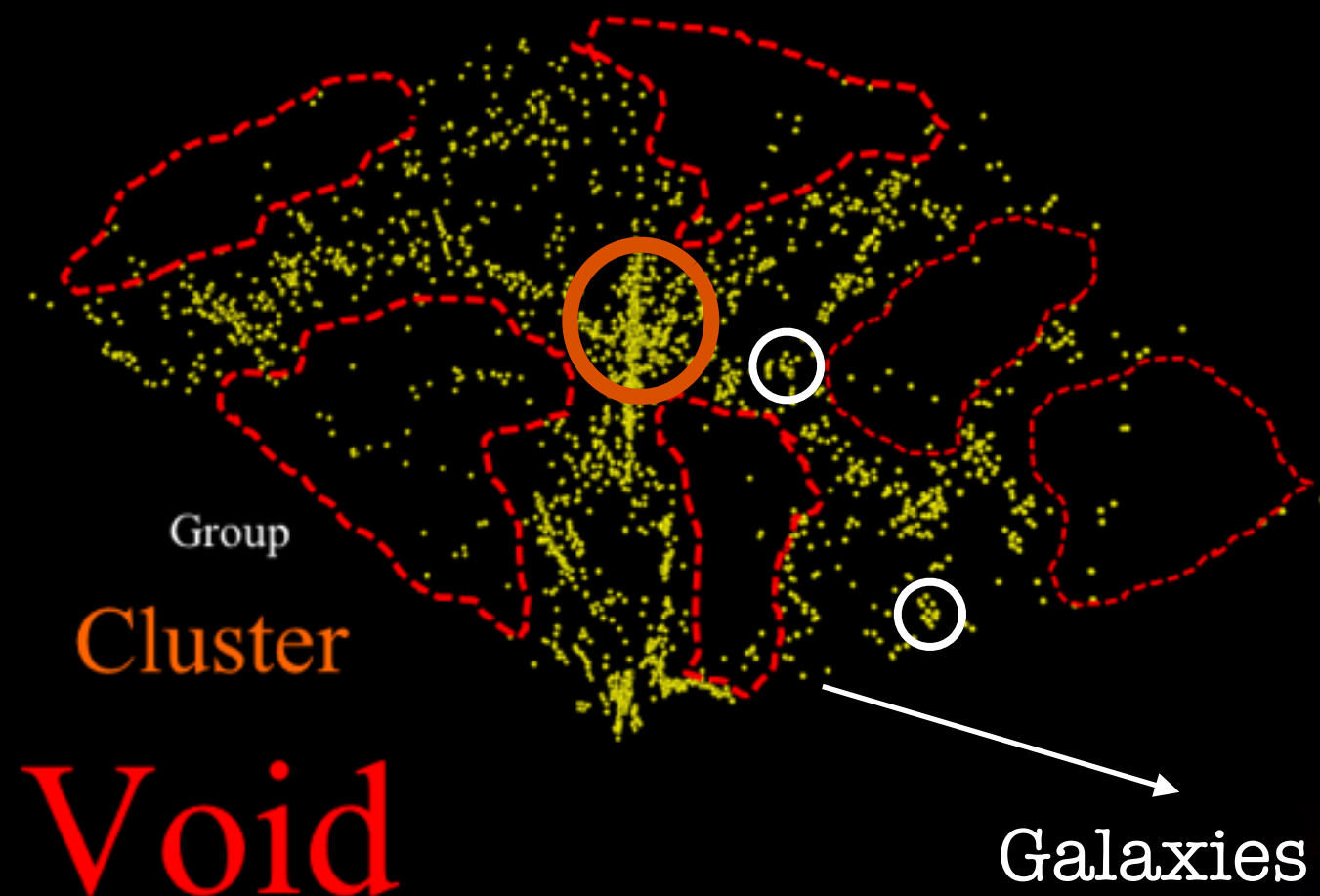
**F. Habibi,  
M. Moniez, R. Ansari & JE. Campagne**

LSST-Webinaire  
June 2017

# Cosmological surveys

**All sky surveys → cosmic structures**

**Deep surveys → structures formation & evolution**



To know about the nature of  
Dark Matter & Dark Energy



# Object classification

- **Cosmic structures contain galaxies.**
- **Images taken by surveys include QSOs and foreground stars in addition to galaxies.**
- **How to separate point-like sources from galaxies?**



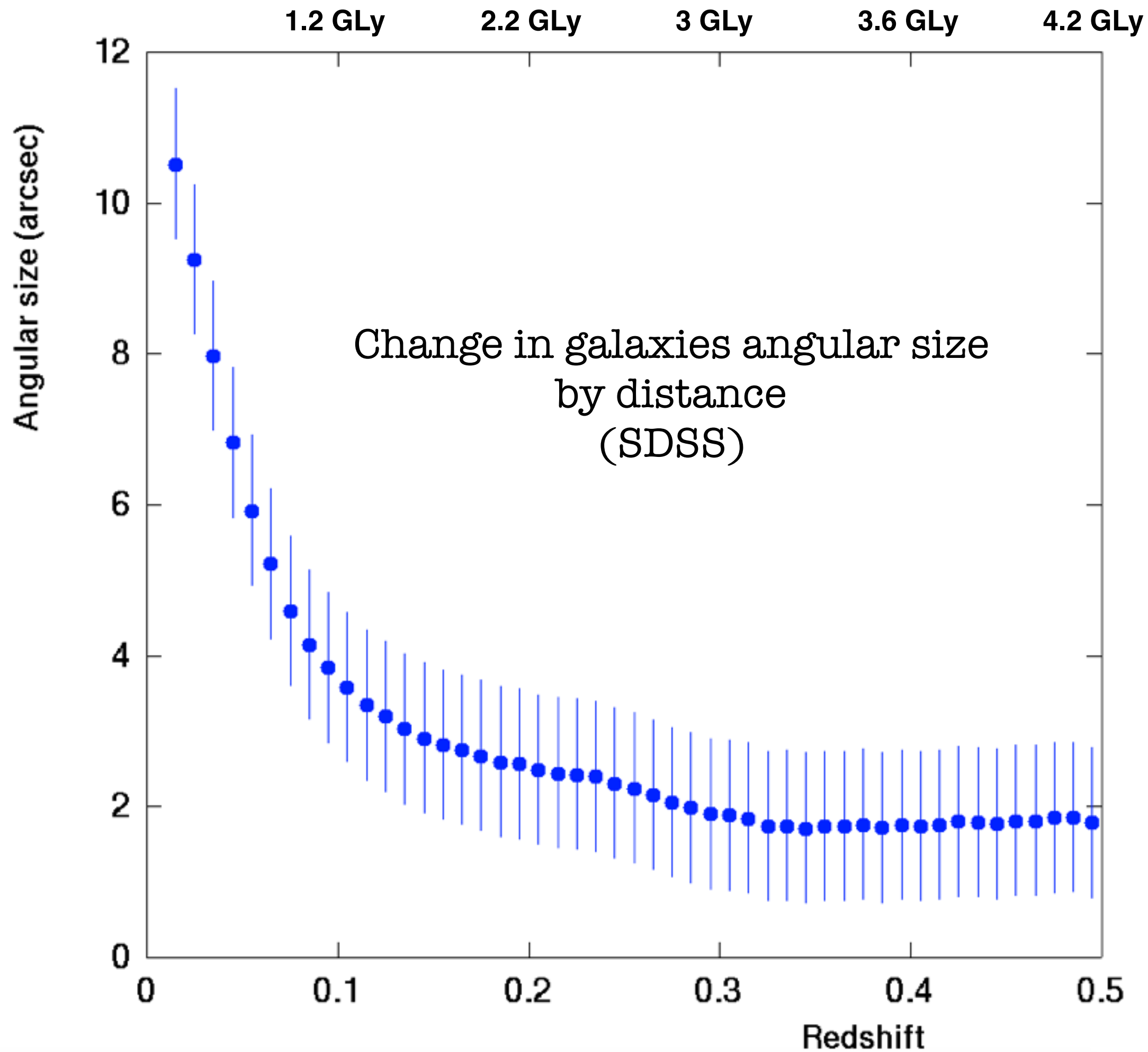
# Nearby galaxies

luminosity spread on CCD:

stars  $\sim 1$  arcsec

galaxies  $\sim 10$  arcsec

full moon  $\sim 1800$  arcsec





# Far/faint galaxies

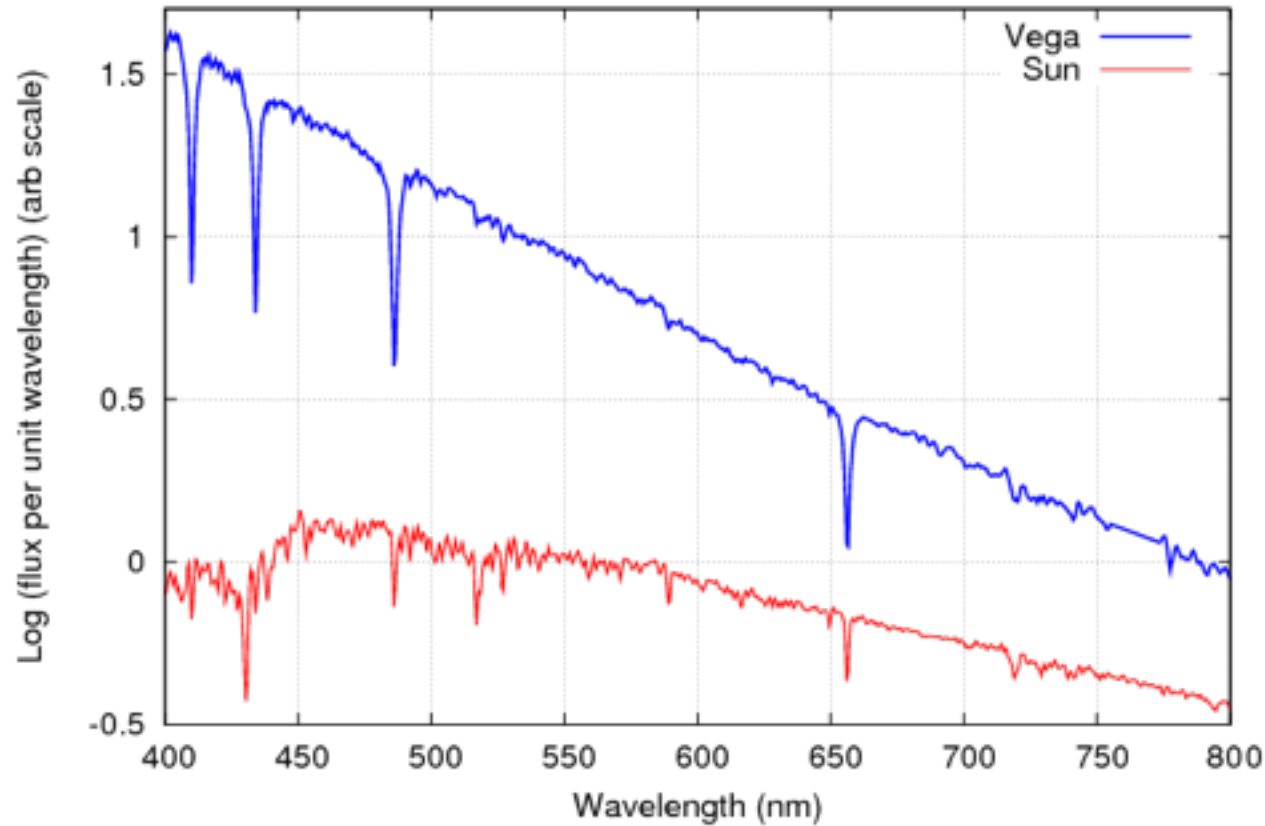
**luminosity spread on CCD:**

**stars  $\sim$  1 arcsec**

**galaxies  $\sim$  1 arcsec**

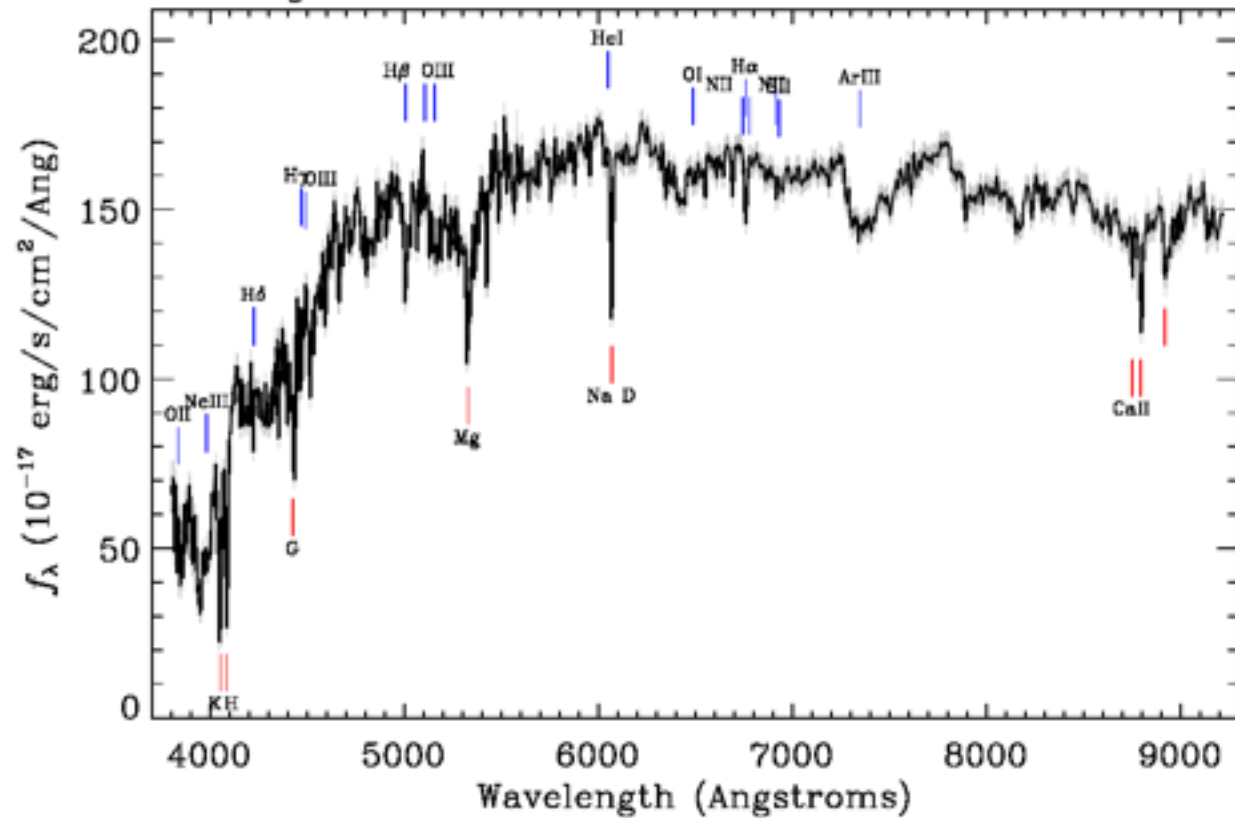


### Spectra of two stars

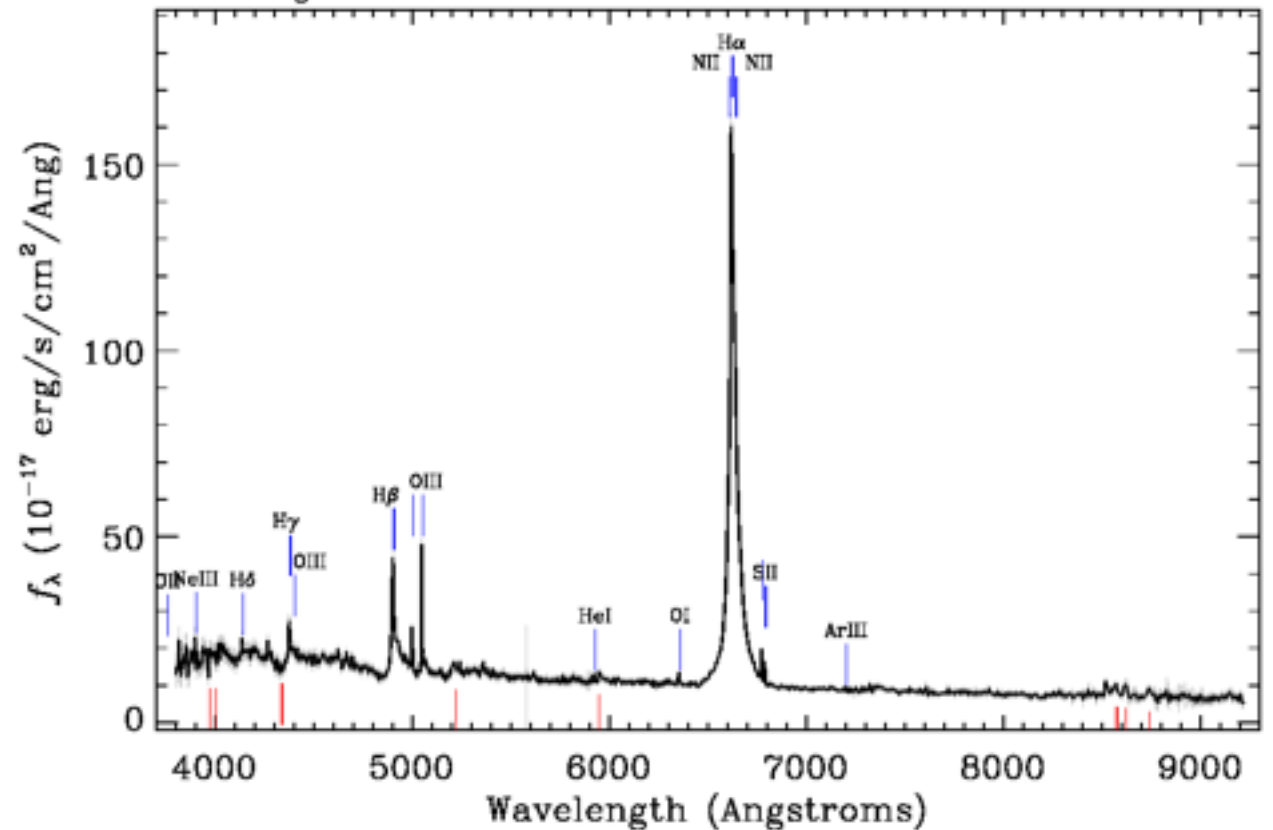


SEDs  
can separate  
the three objects

Survey: sdss Program: legacy Target: GALAXY\_RED GALAXY  
RA=242.84830, Dec=52.45478, Plate=623, Fiber=433, MJD=52051  
z=0.02942±0.00001 Class=GALAXY  
No warnings.



Survey: sdss Program: legacy Target: QSO\_CAP  
RA=173.34990, Dec=55.07108, Plate=1014, Fiber=463, MJD=52707  
z=0.00906±0.00009 Class=QSO BROADLINE  
No warnings.



# Aim

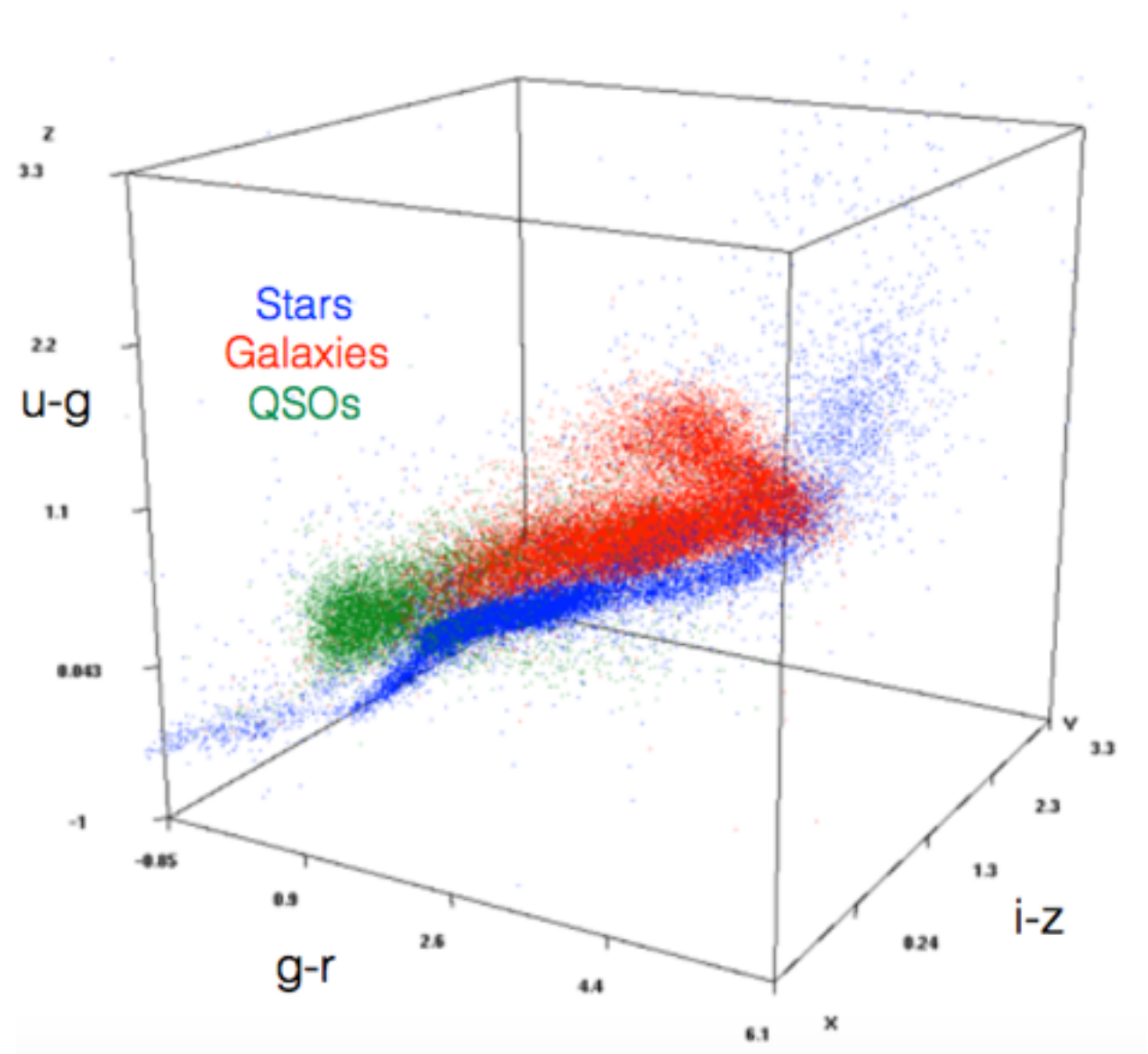
To separate galaxies from stars and QSOs,  
in the lack of spectroscopic data.



# How?

Including all possible photometric information

**Colour indices:**  
**proper “features”**  
**for supervised**  
**perceptrons**

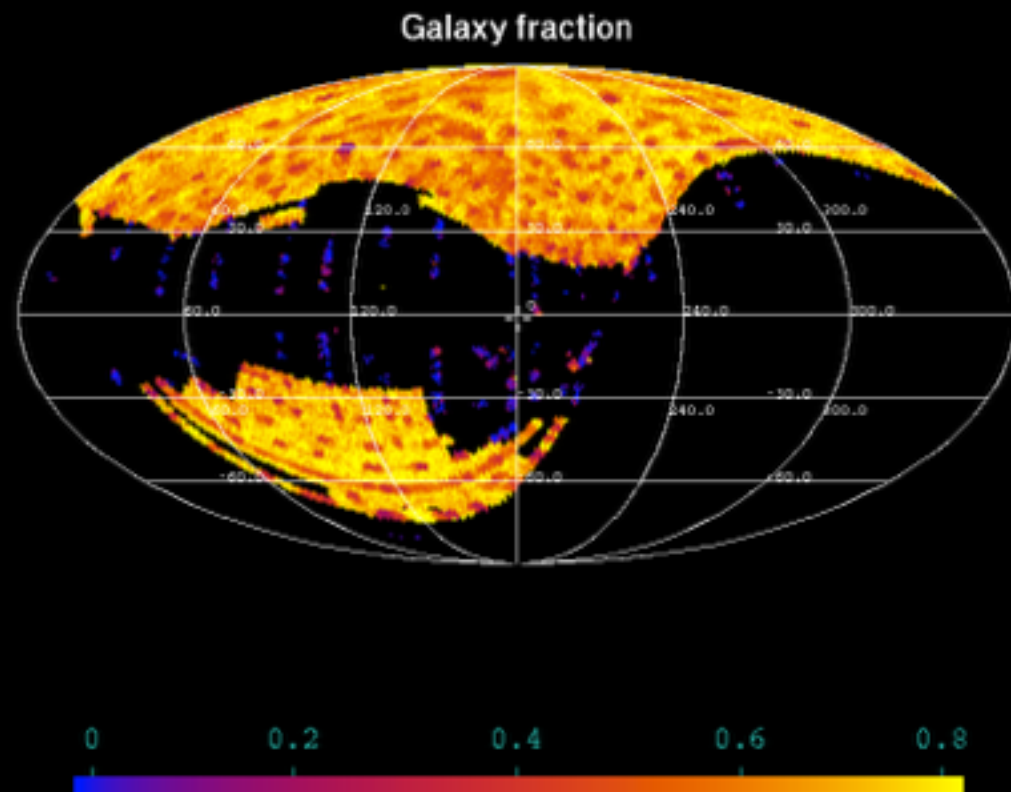


Using automatic classification for big number of data

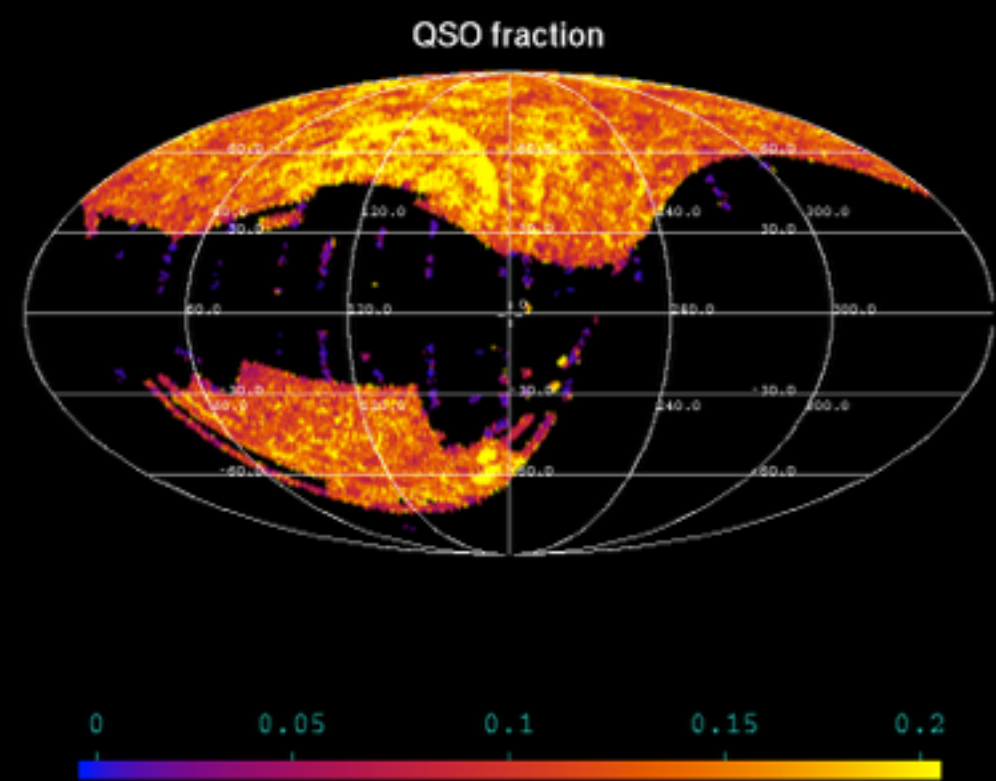
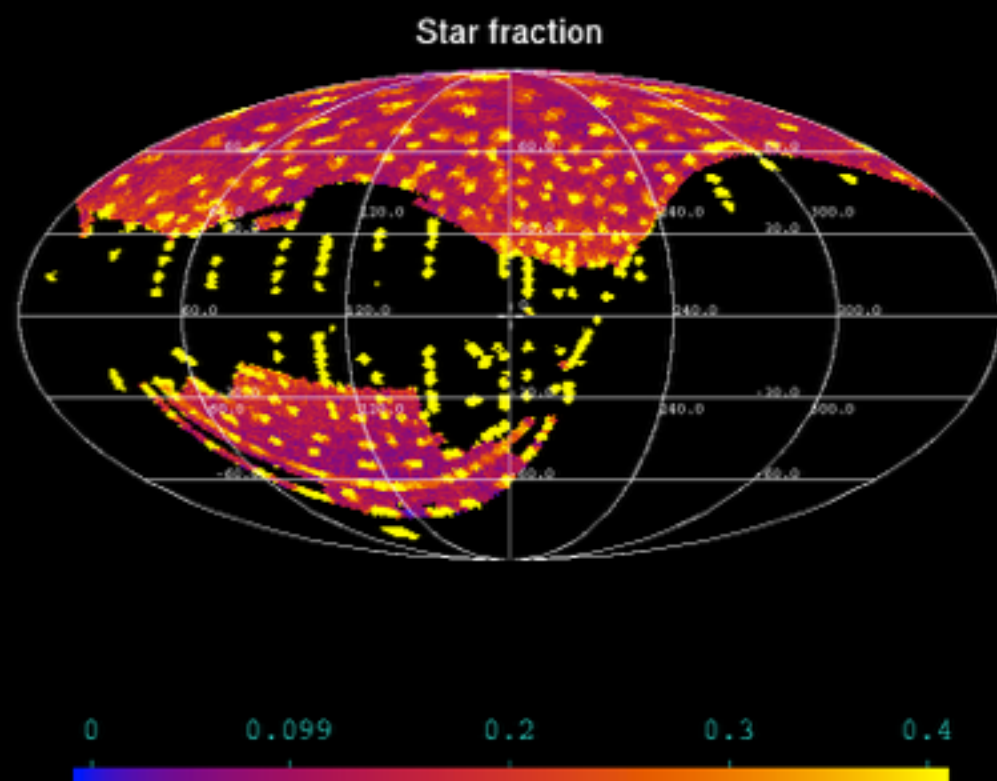
# SDSS DR12 data

## Different sub-surveys

Objects selected with both  
photometric and spectroscopic  
data available



Stars	Galaxies	QSOs	Total
928,464	2,484,161	566,475	3,979,100
23%	62%	15%	





# SDSS DR12 photometry

- PSF magnitude
- **Model magnitude**: de Vaucouleurs / exponential profile

$$I(r) = I_0 \exp\{-7.67 [(r/r_e)^{1/4}]\}$$

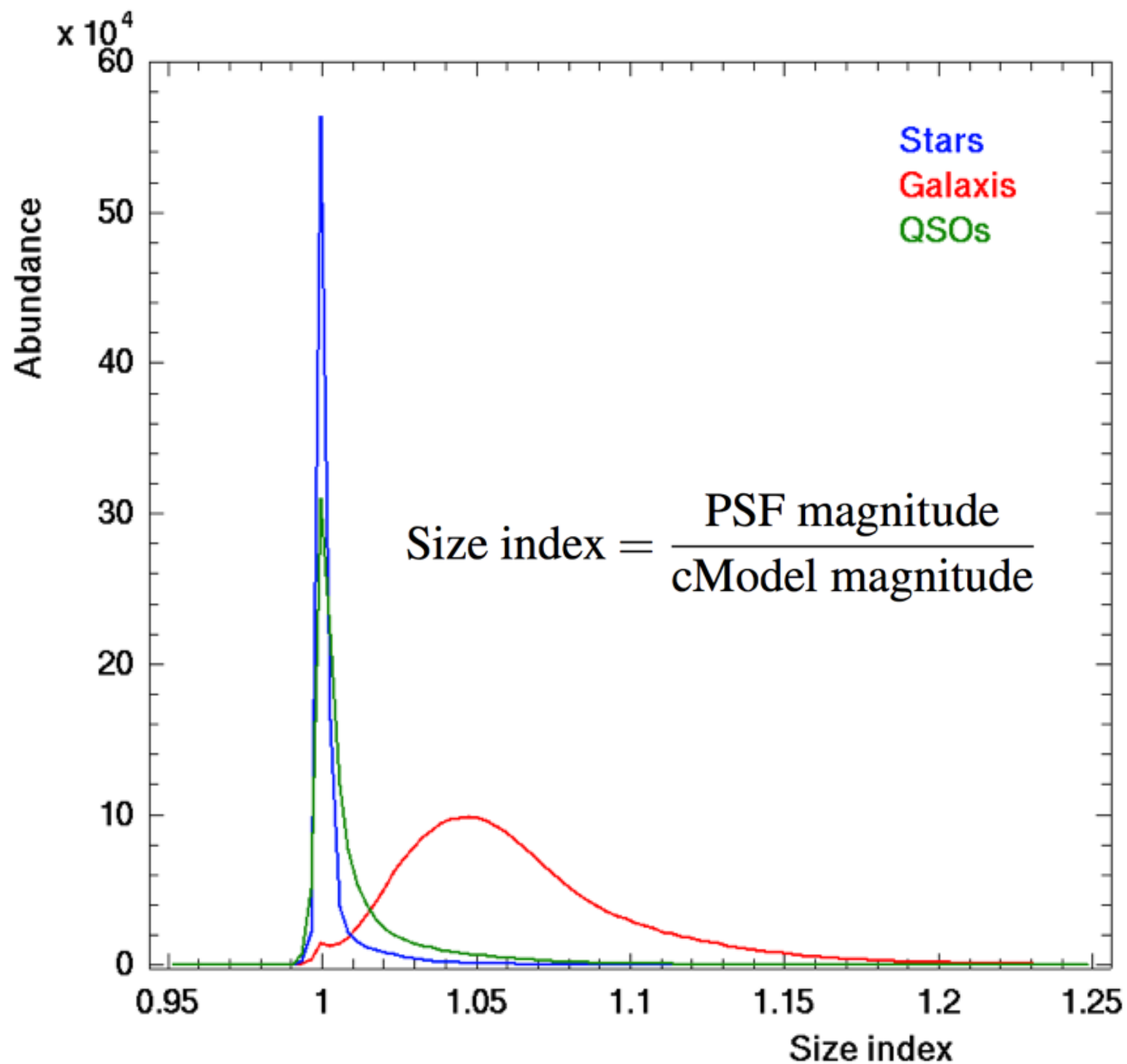
$$I(r) = I_0 \exp(-1.68r/r_e)$$

- **Composite model magnitude**:

$$F_{\text{composite}} = \text{fracDev} F_{\text{dev}} + (1 - \text{fracDev}) F_{\text{exp}}$$

$$\text{Size index} = \frac{\text{PSF magnitude}}{\text{cModel magnitude}}$$

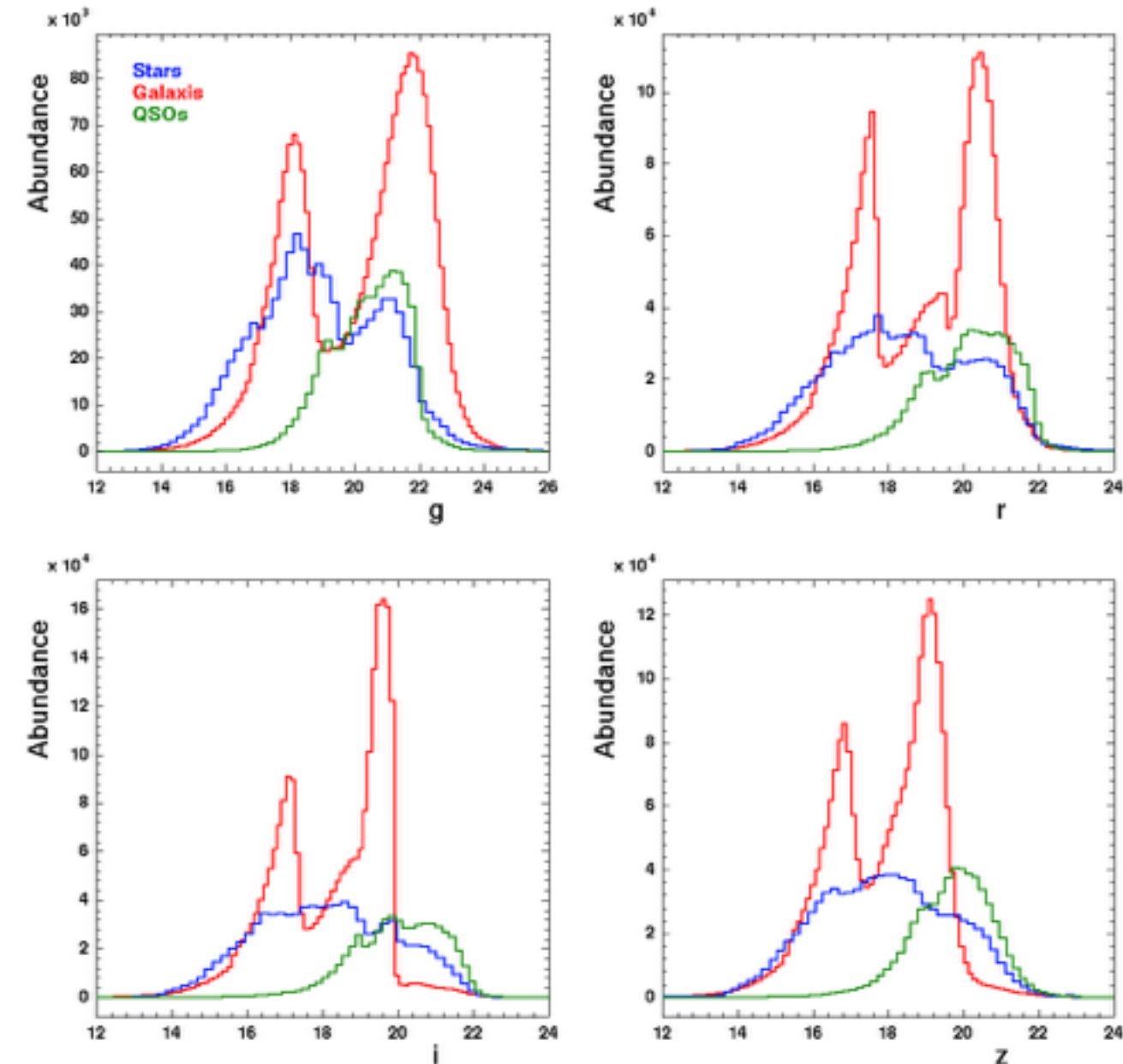
# SDSS DR12 data





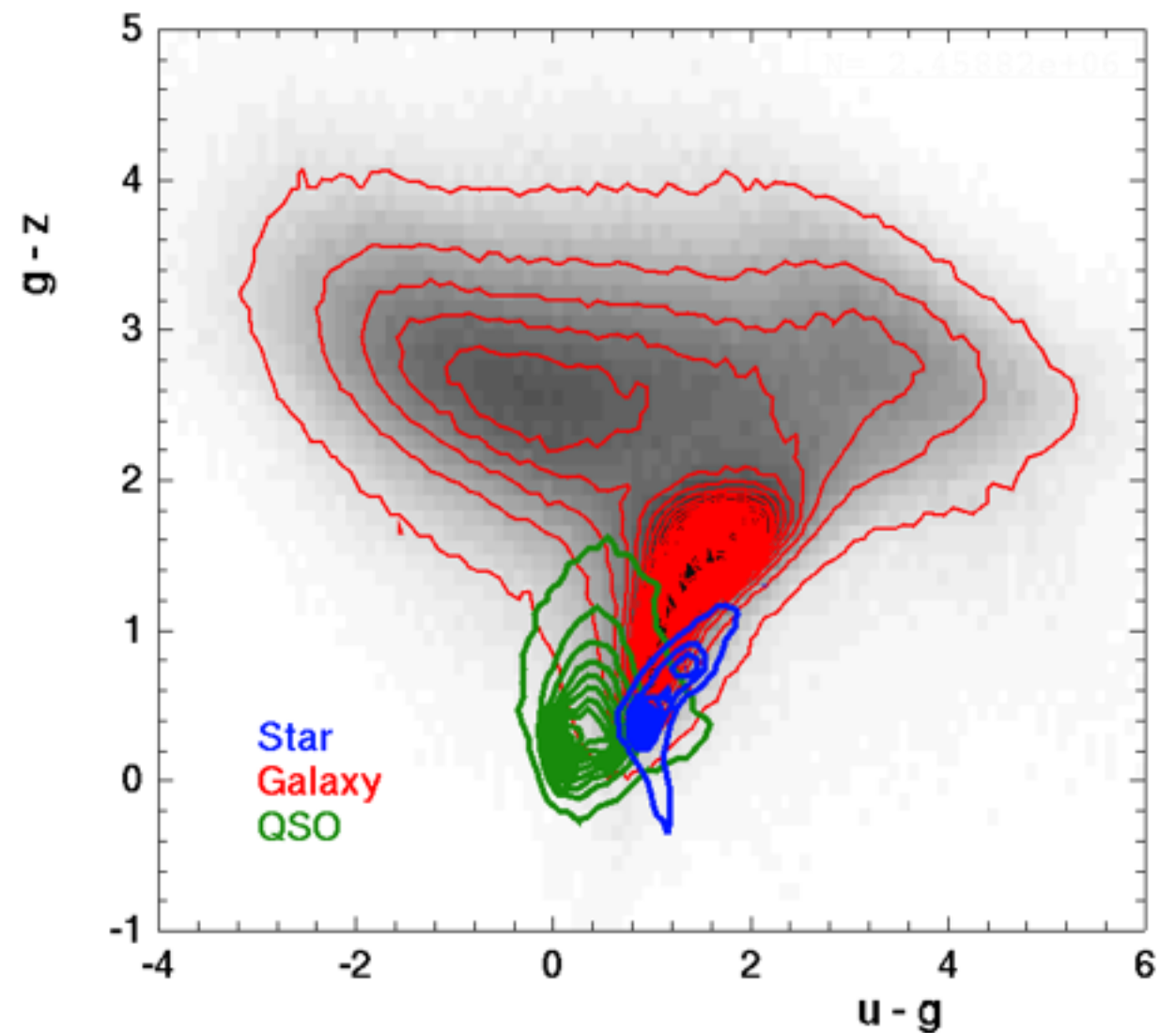
# SDSS DR12 data

## Magnitude distributions



Magnitudes:  
insufficient to separate the objects

## Colour-colour diagram



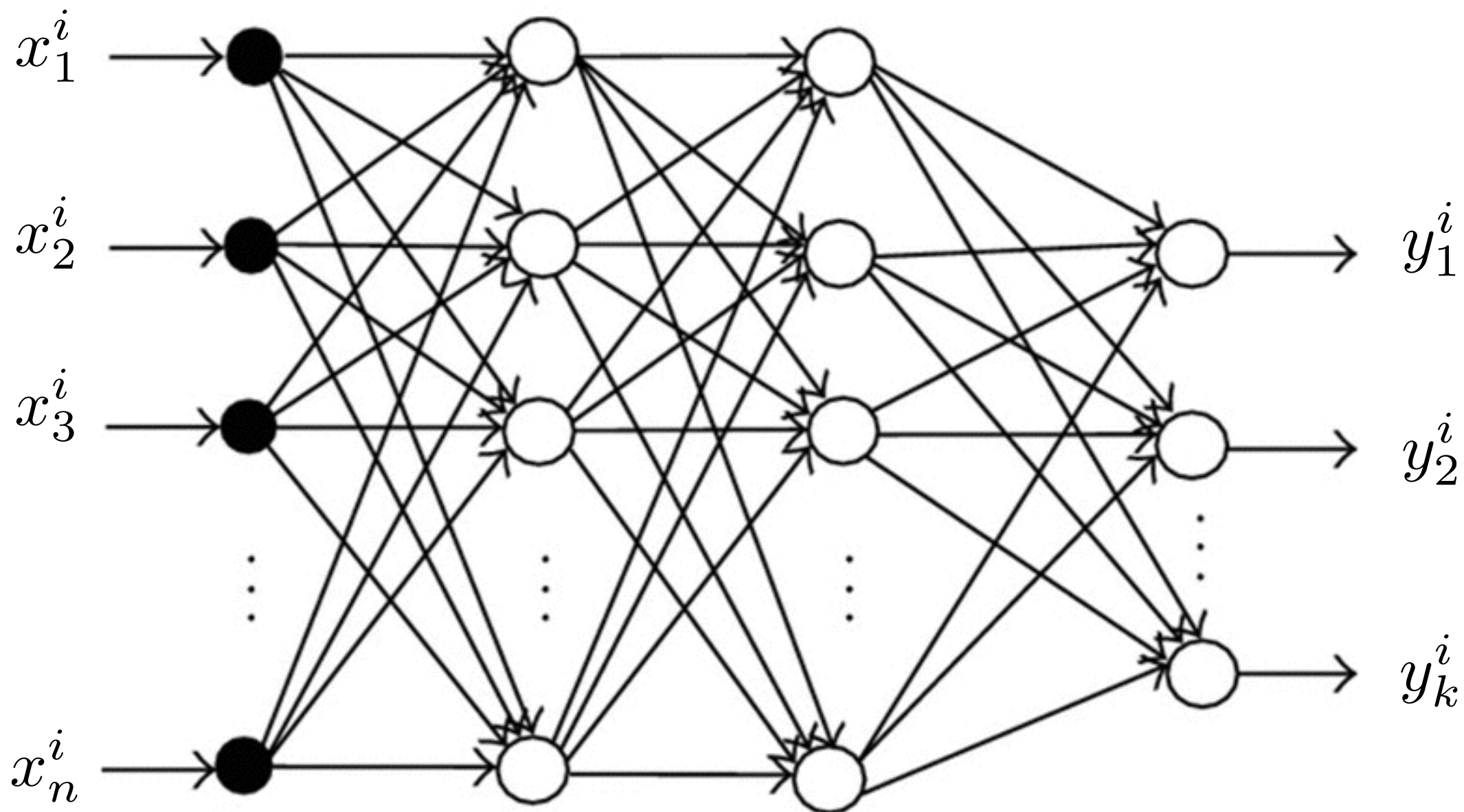
galaxies are redder  
in average

# Multi layer perceptron

$m$  objects in the training set,  $1 < i < m$

Each object contains  $n$  features

$k$  number of classes labelled by vector  $y^i$



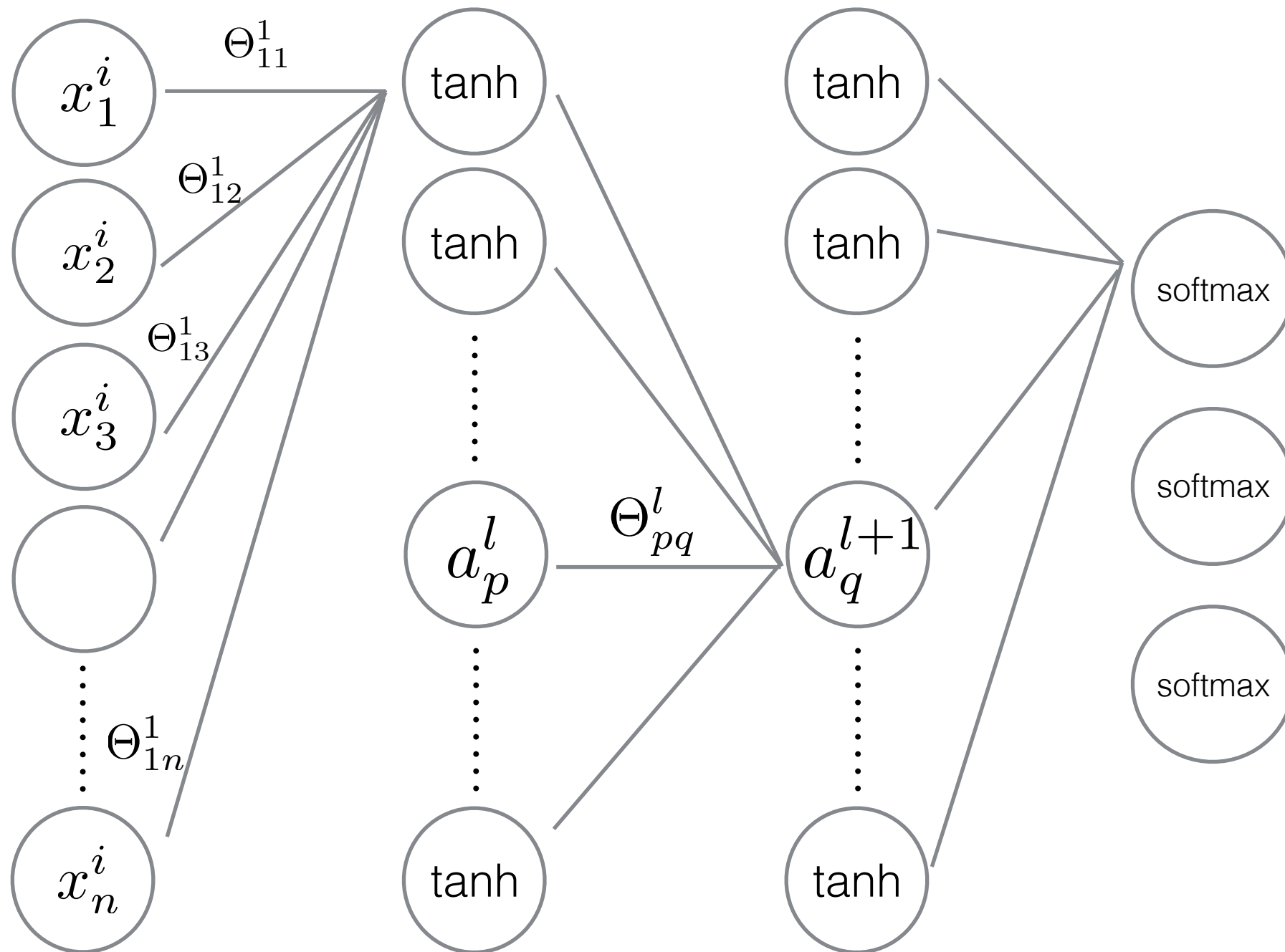
$$J(\Theta) = -\frac{1}{m} \left[ \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log((h_{\Theta}(x^{(i)}))_k) + (1 - y_k^{(i)}) \log(1 - (h_{\Theta}(x^{(i)}))_k) \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\Theta_{j,i}^{(l)})^2$$



# Multi layer perceptron

$$a_q^{l+1}(x^i) = g\left(\sum_{p=0}^{s_l} \Theta_{pq}^l a_p^l(x^i)\right)$$

$g(z)$ : Activation function: Sigmoid, tanh, softmax and etc.



$$J(\Theta) = -\frac{1}{m} \left[ \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log((h_{\Theta}(x^{(i)}))_k) + (1 - y_k^{(i)}) \log(1 - (h_{\Theta}(x^{(i)}))_k) \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\Theta_{j,i}^{(l)})^2$$

# Activation

## Sigmoid function

$$h(z) = \frac{1}{1 + e^{-z}}$$

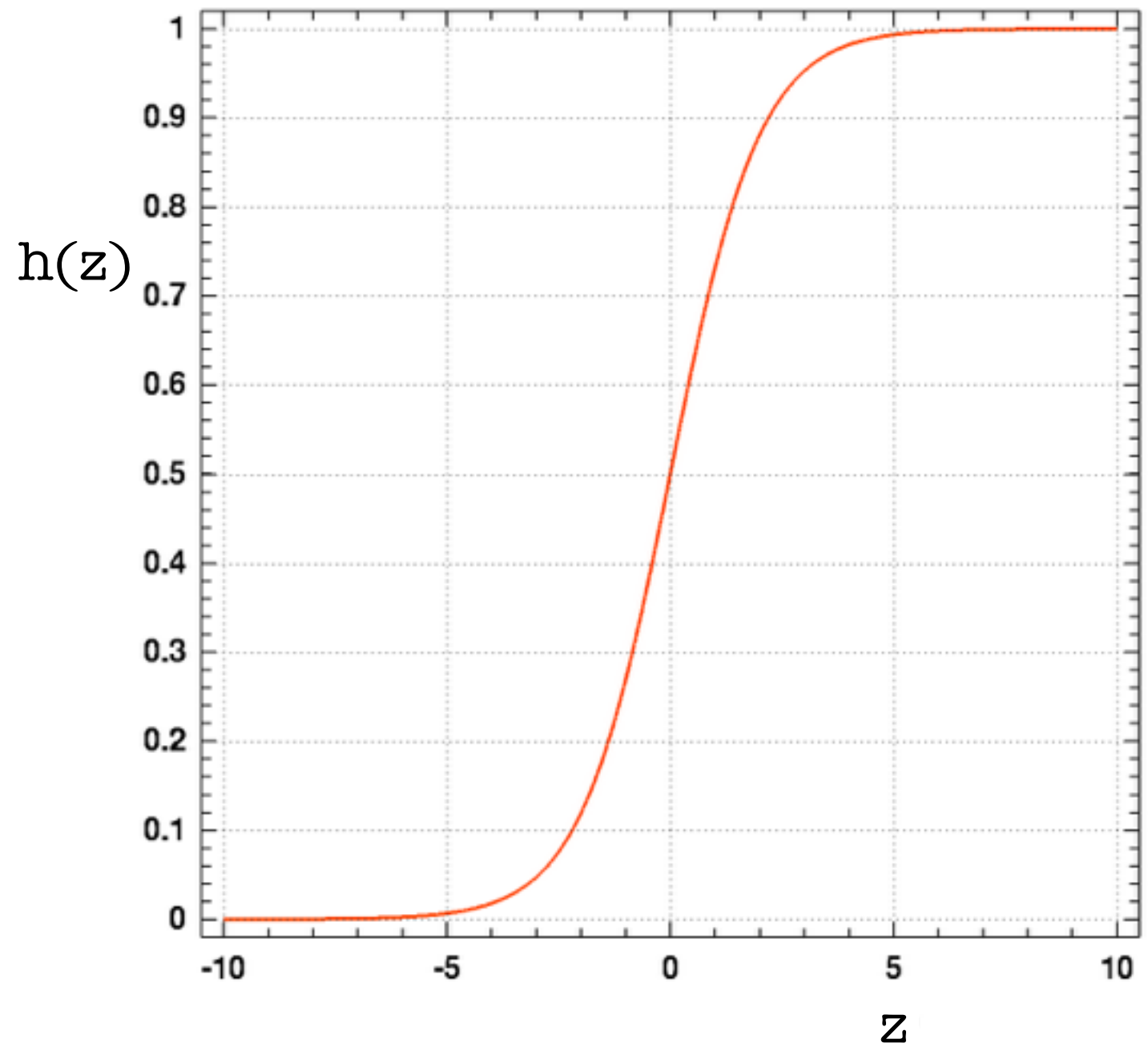
$$z = \theta_0 + \theta_1 x_1 + \dots + \theta_n x_n$$

Border hyper surface:

$$h(z = 0) = 0.5$$

A single-minimum  
cost function:

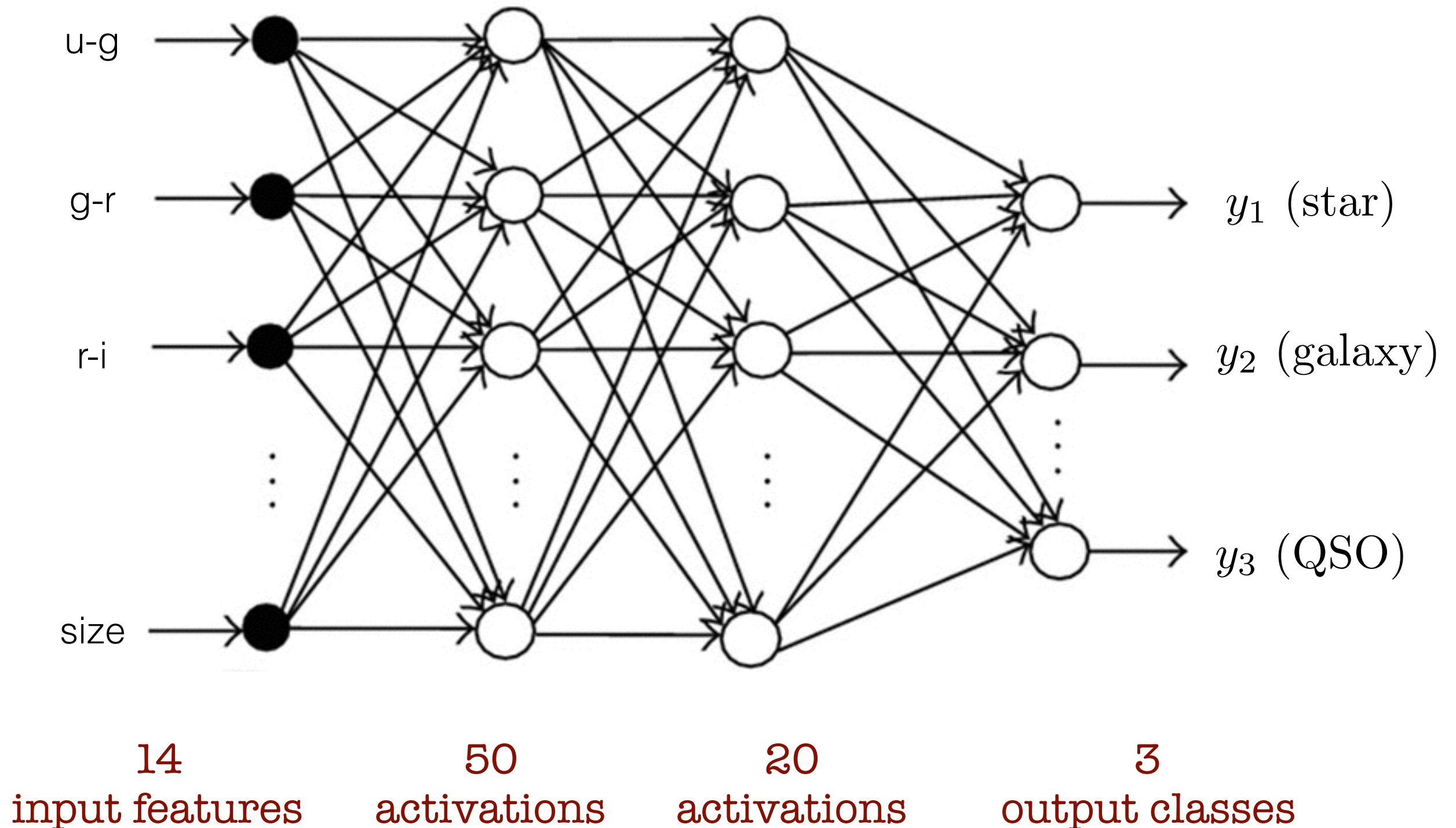
$$J = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(h_{\theta}(\vec{x}^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(\vec{x}^{(i)}))]$$



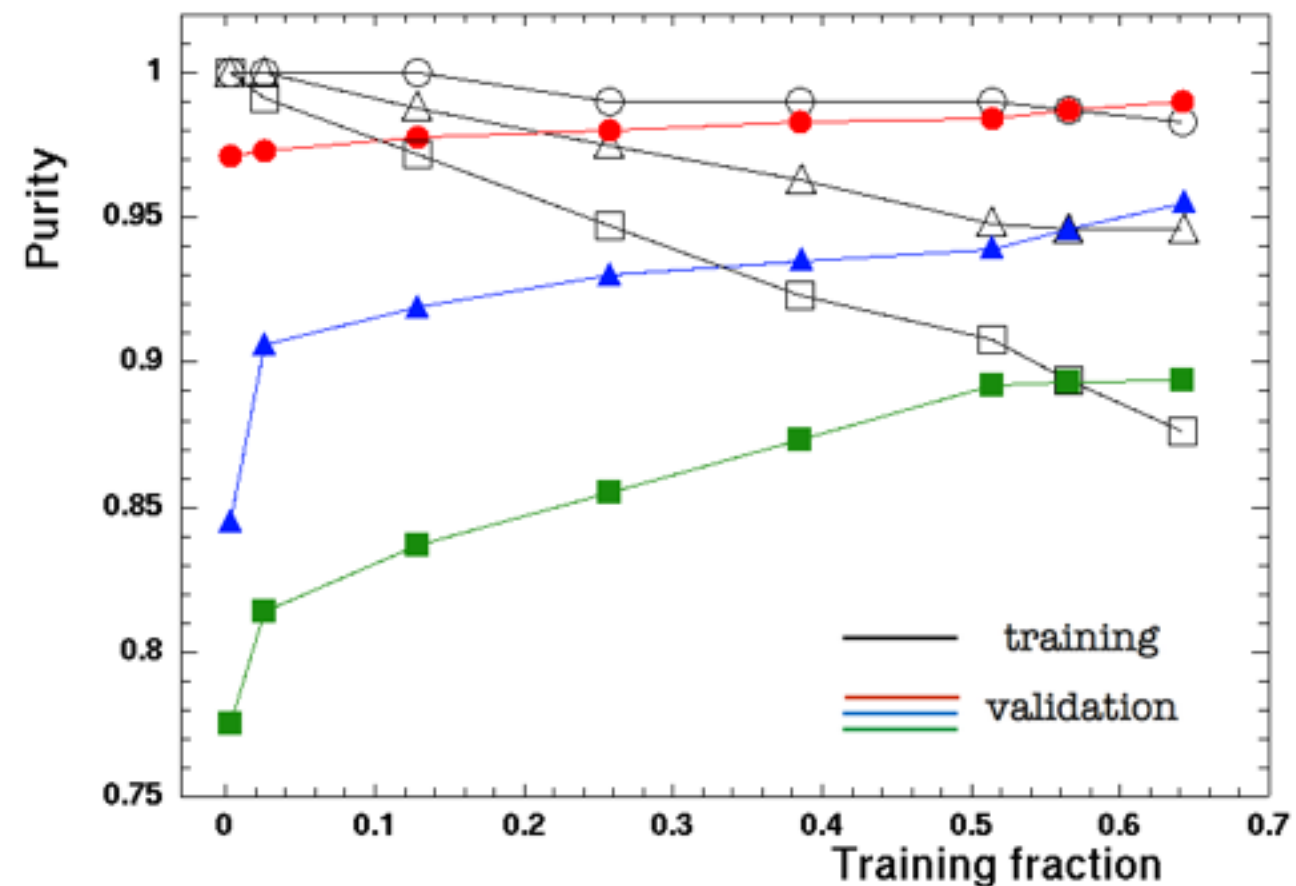
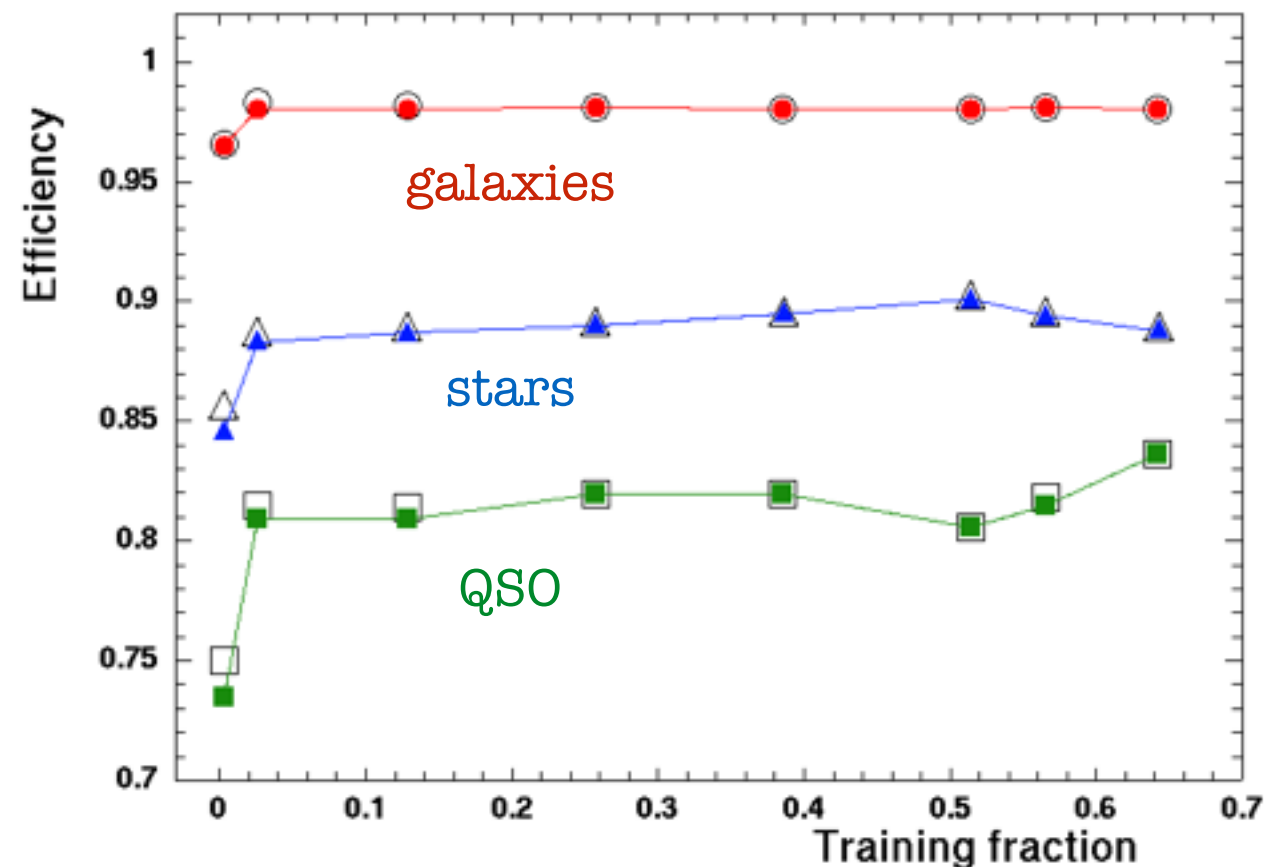


# Multi layer perceptron

Dense Neural Net with 2 hidden layers



# Training the MLP



$$\text{efficiency}_i = \frac{n_{i \rightarrow i}}{n_i}$$

$$\text{purity}_i = \frac{n_{i \rightarrow i}}{n_{i \rightarrow i} + \sum_{j \neq i} n_{j \rightarrow i}}$$

$i$  : galaxy, star or QSO

Comparing  
training and validation sets  
to find optimum number  
of objects for training



# Training the MLP

## Efficiency and purity of the classification

size index included

	Star	Galaxy	QSO	Total
Efficiency	89.4%	98.1%	81.5%	94%
Purity	94.6%	98.7%	89.3%	-

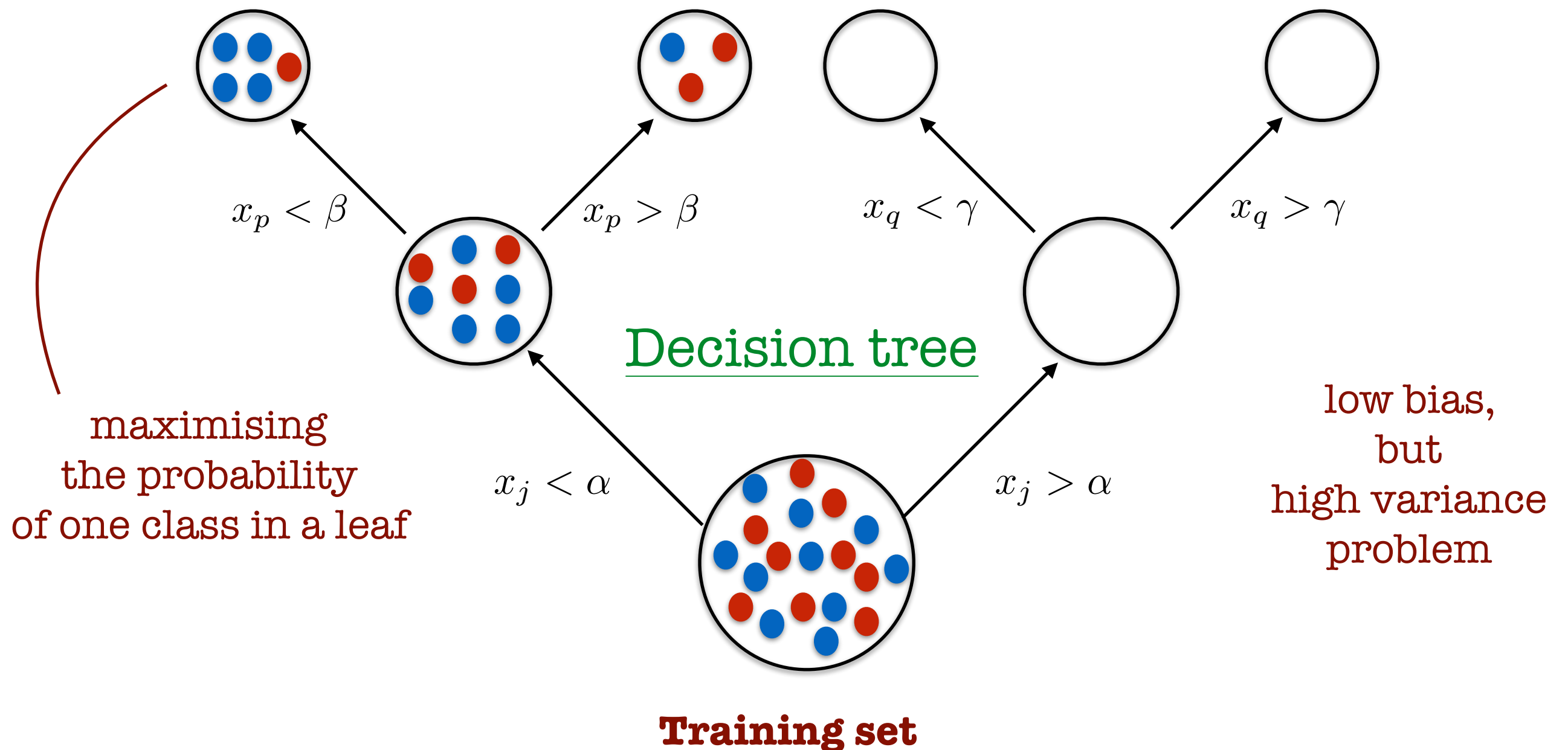
$$\text{Size index} = \frac{\text{PSF magnitude}}{\text{cModel magnitude}}$$

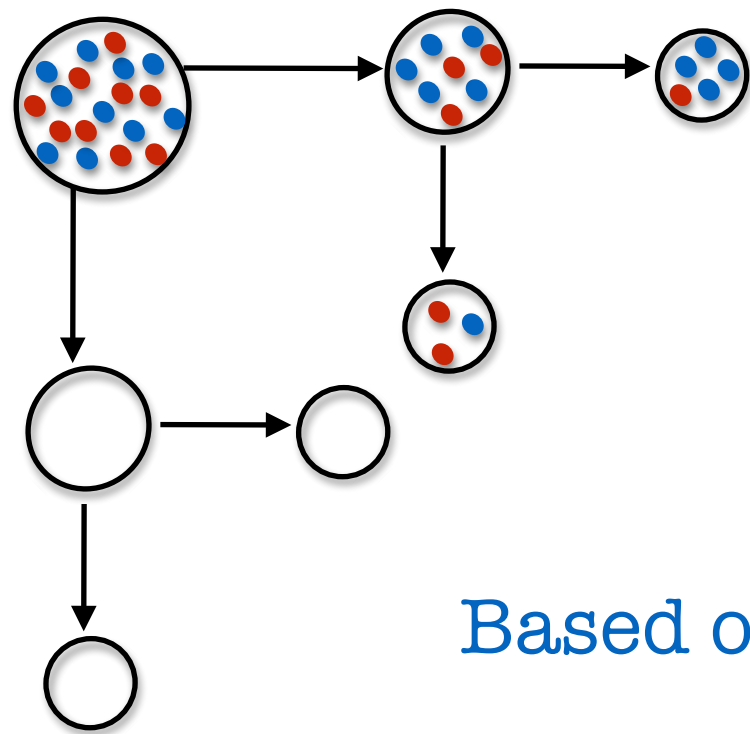
no size index

	Star	Galaxy	QSO	Total
Efficiency	86.6%	97.5%	78.5%	92%
Purity	93.0%	97.7%	88.1%	-

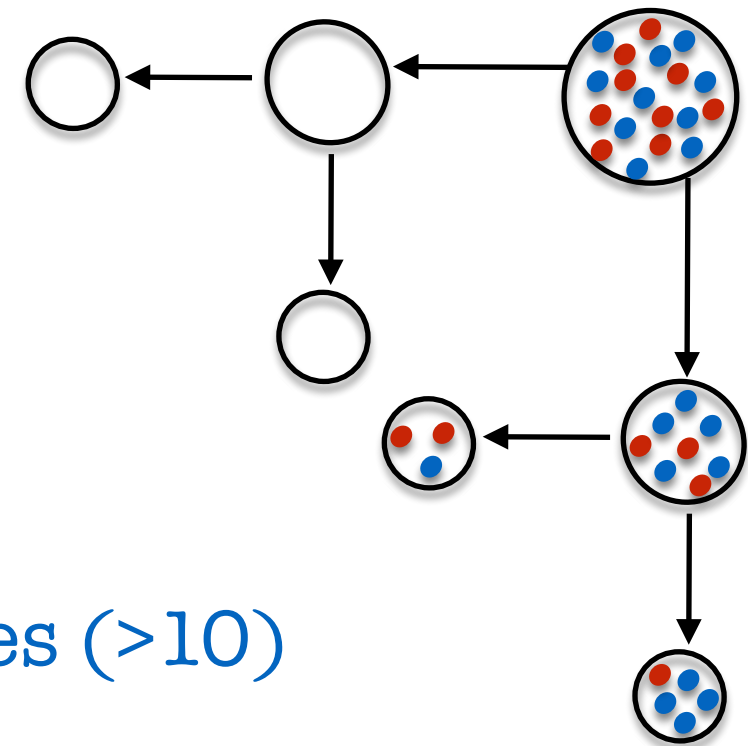
# Random forest

constructing a classification model through feature filtering





# Random forest

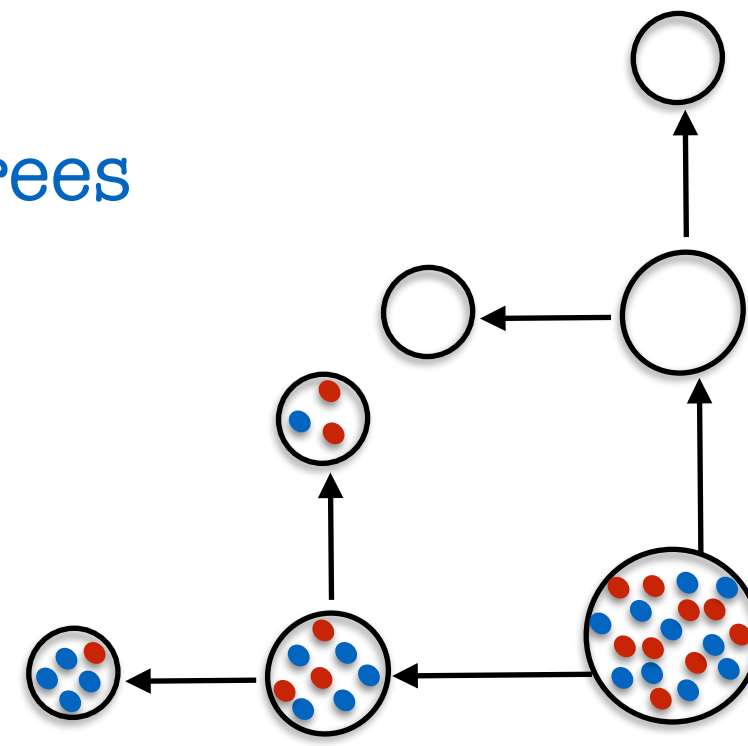
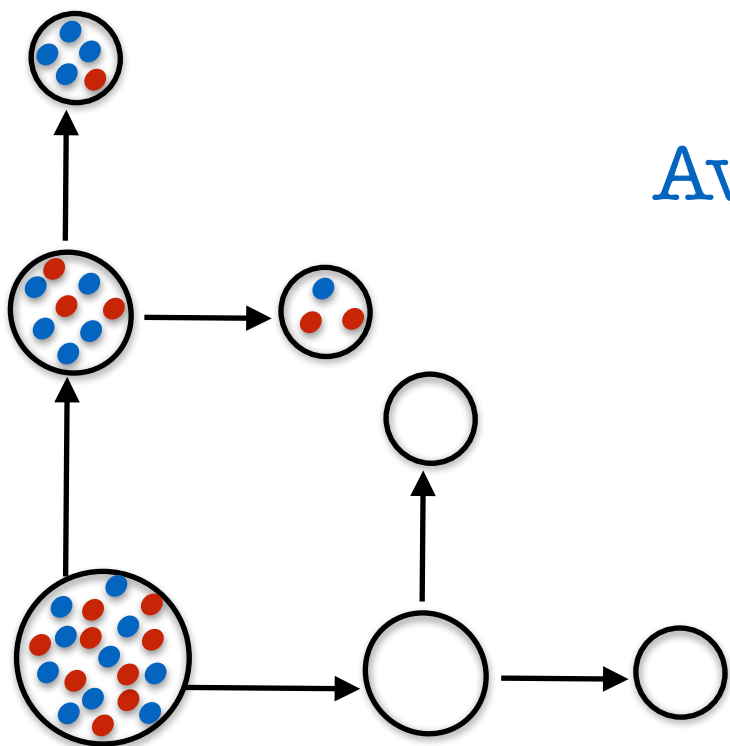


Based on large number of decision trees (>10)

Generating different samples of the training set through bootstrapping (sampling with replacement)

Feature (random) bagging at conjunctions

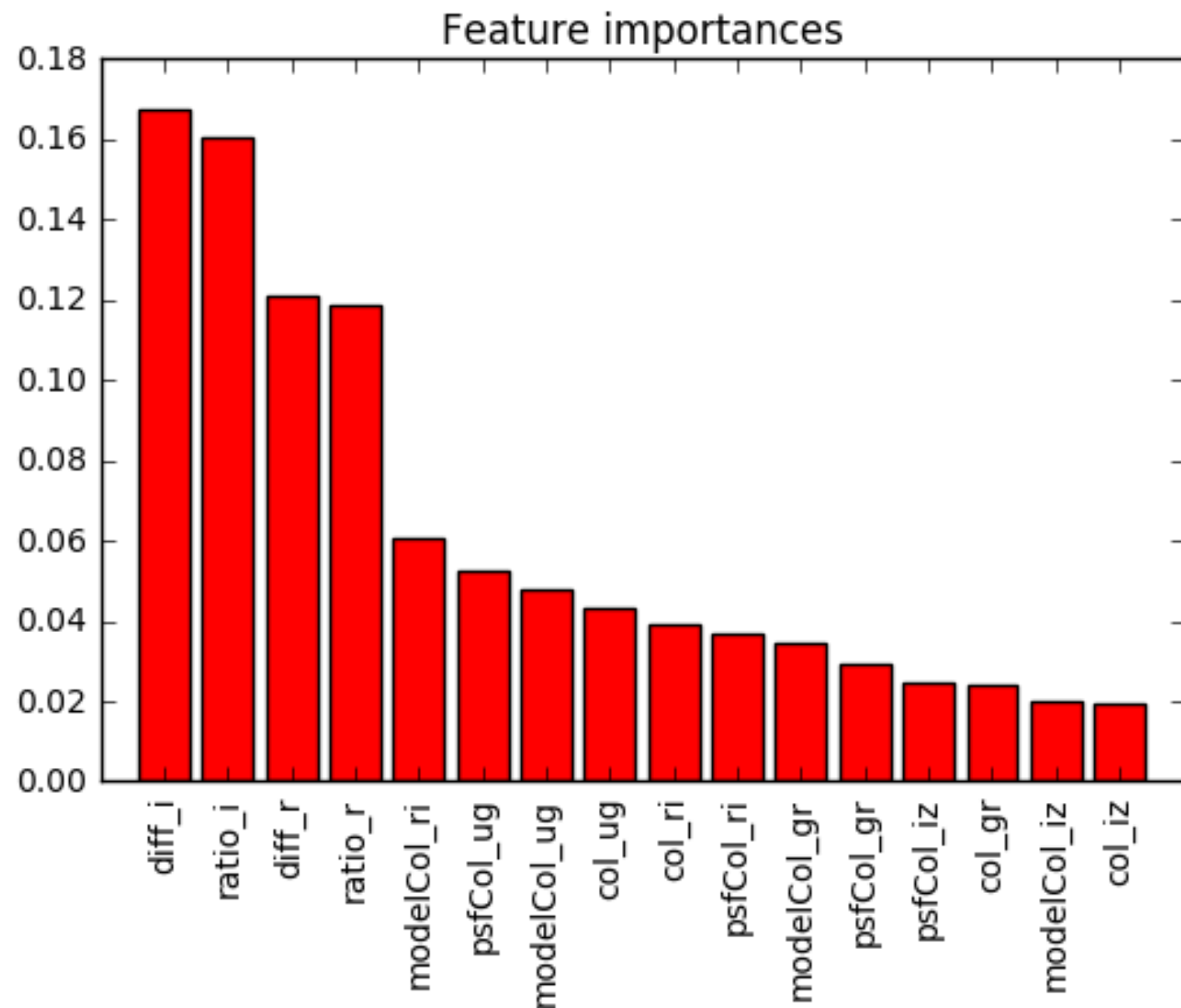
Average over predictions of all trees



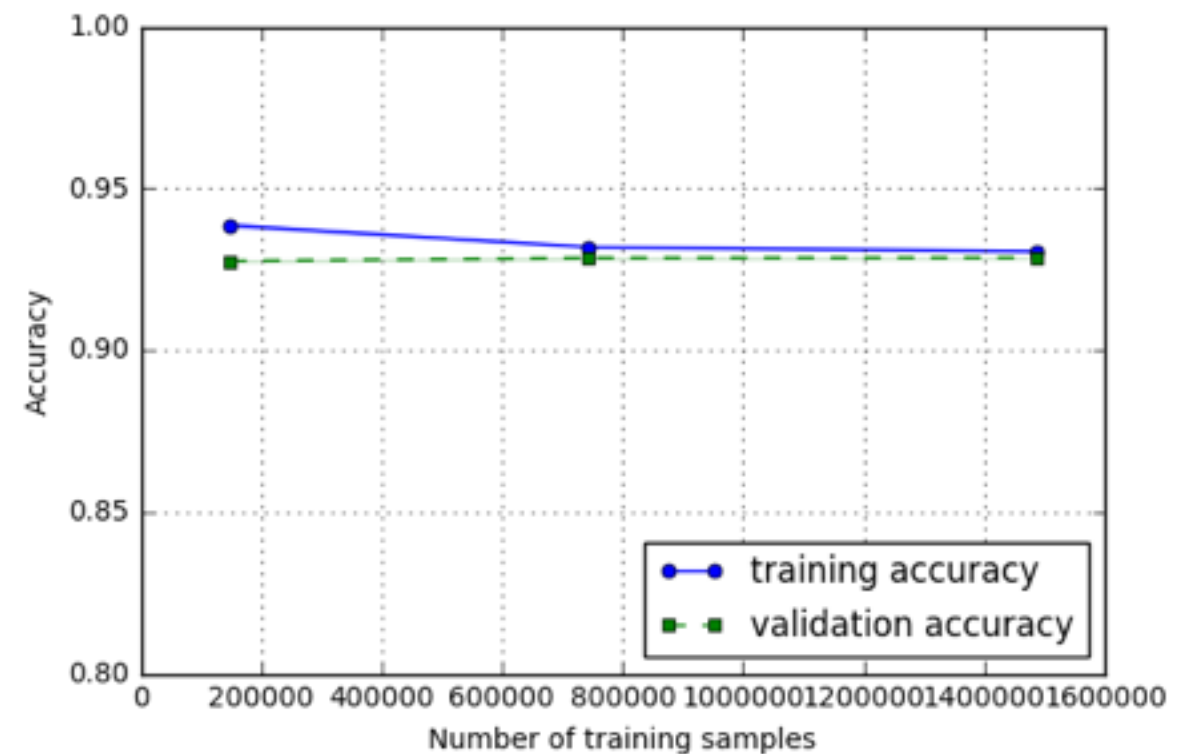


# Random forest

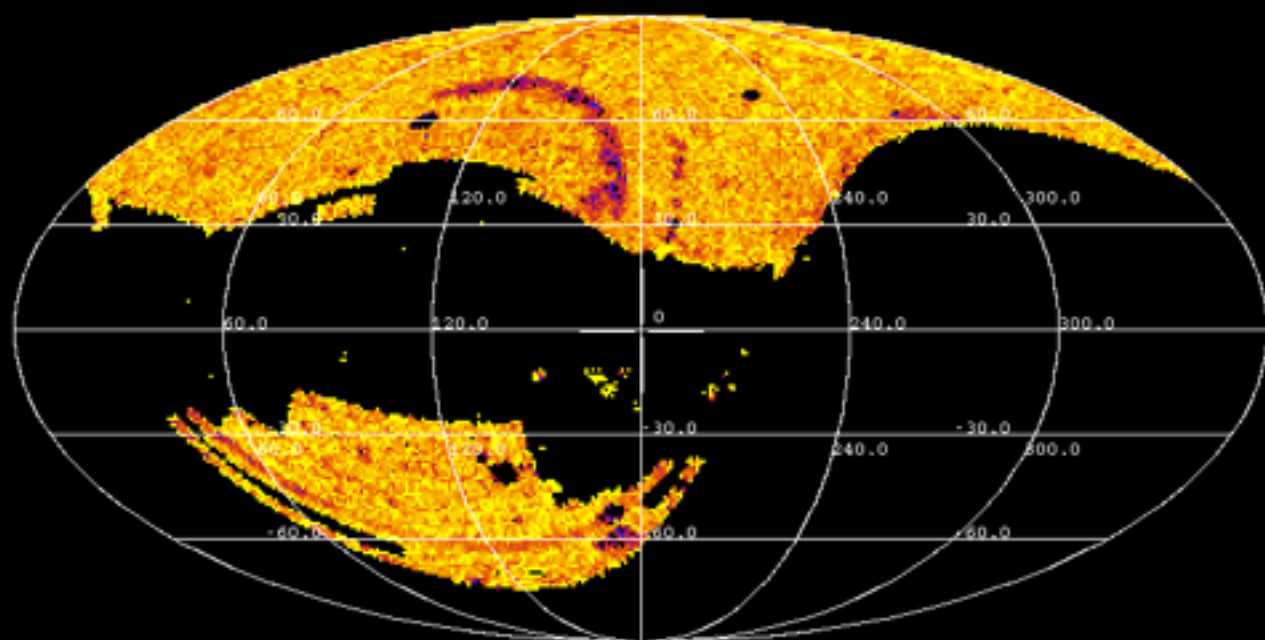
Efficiency	Star	Galaxy	QSO	Total
RF	86.9%	98.0%	80.2%	93%
NN	89.4%	98.1%	81.5%	94%



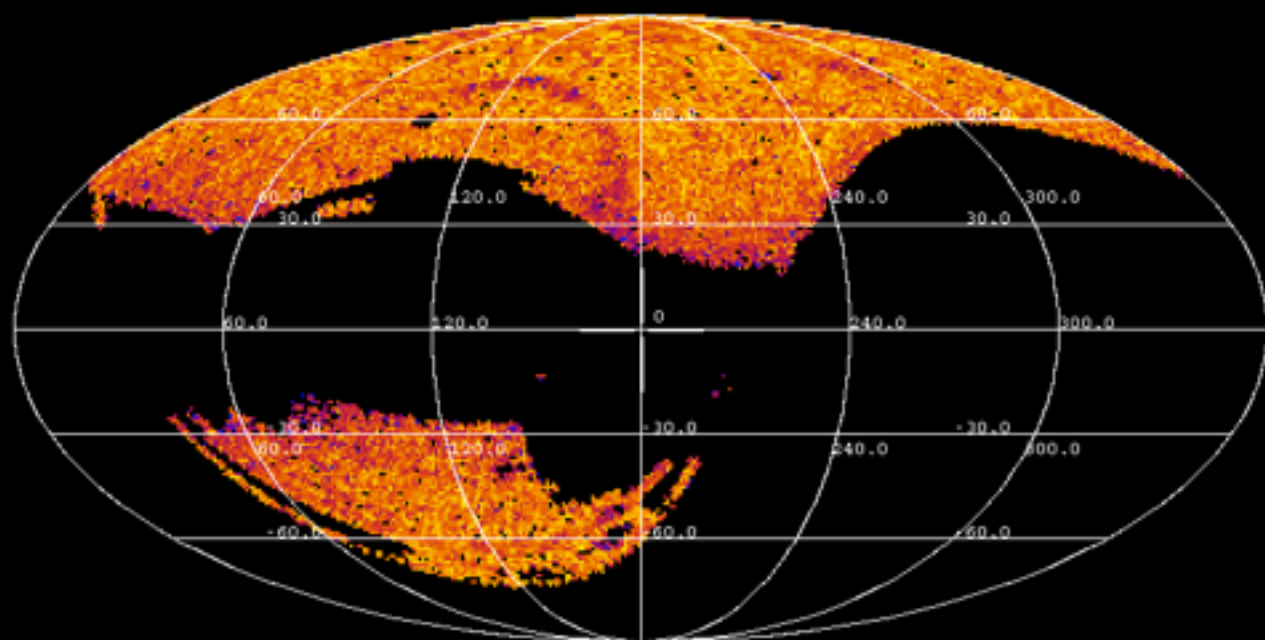
Size index is the most discriminative feature for RF



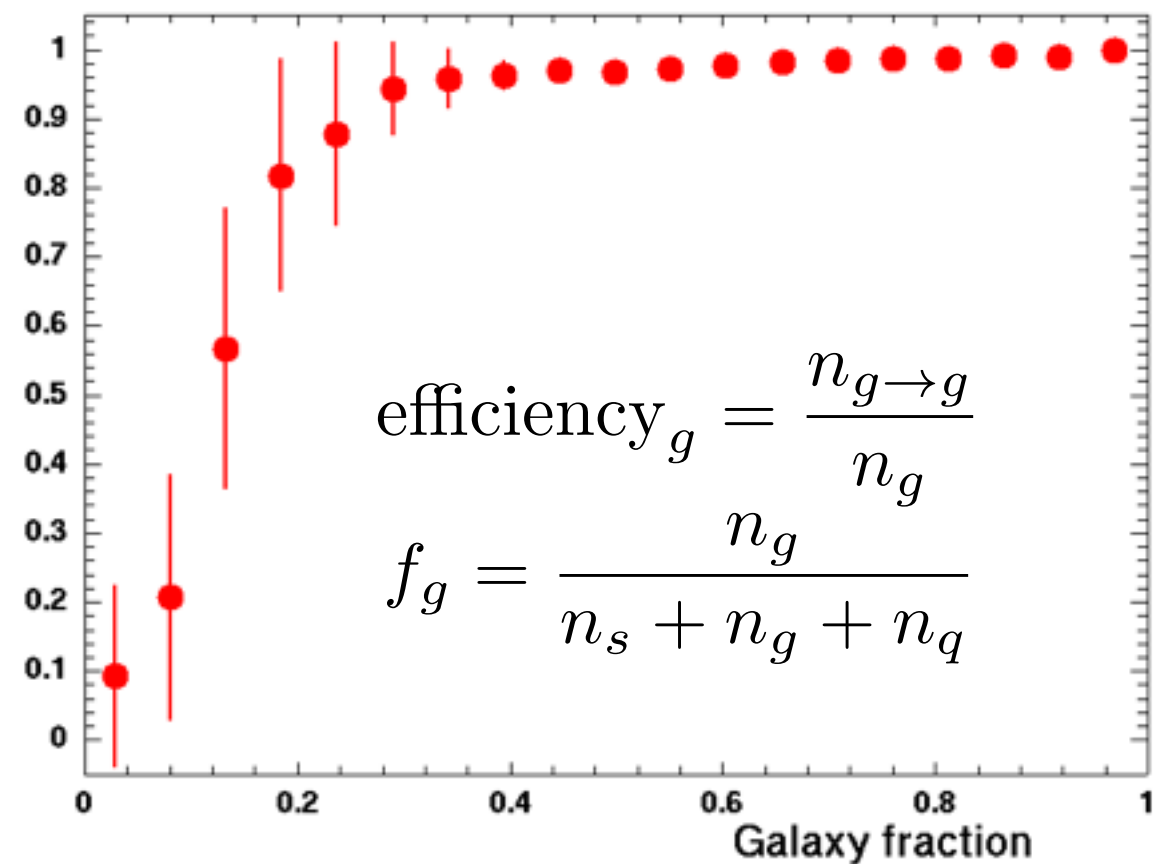
Galaxy efficiency



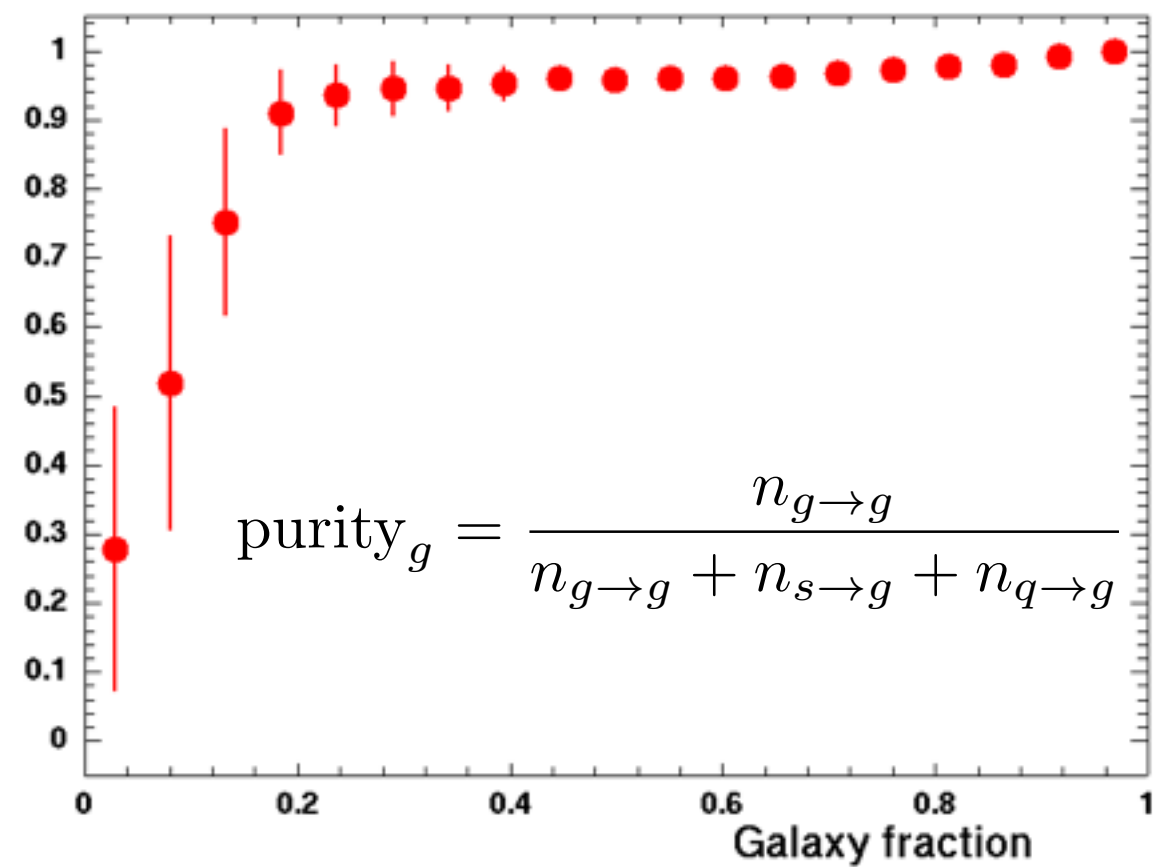
Galaxy purity

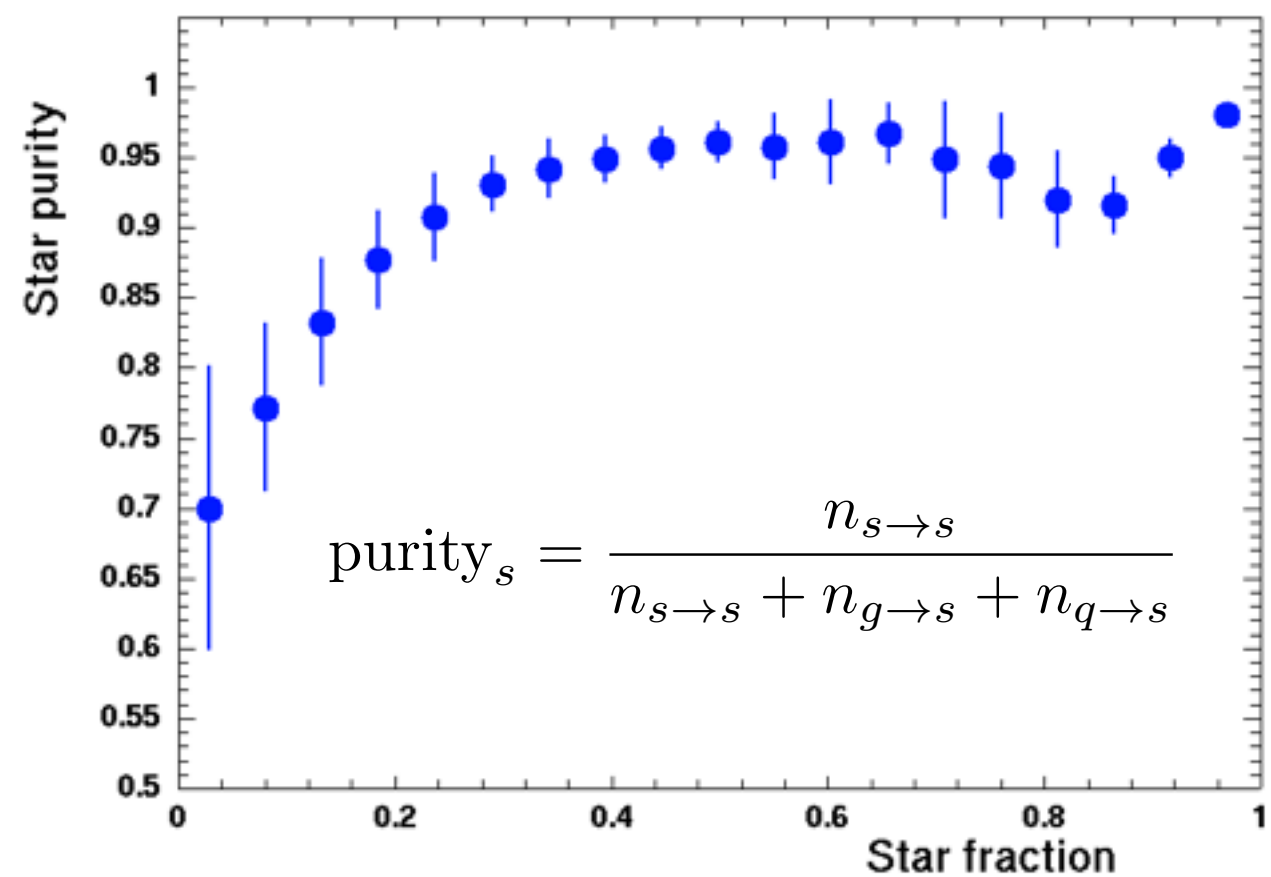
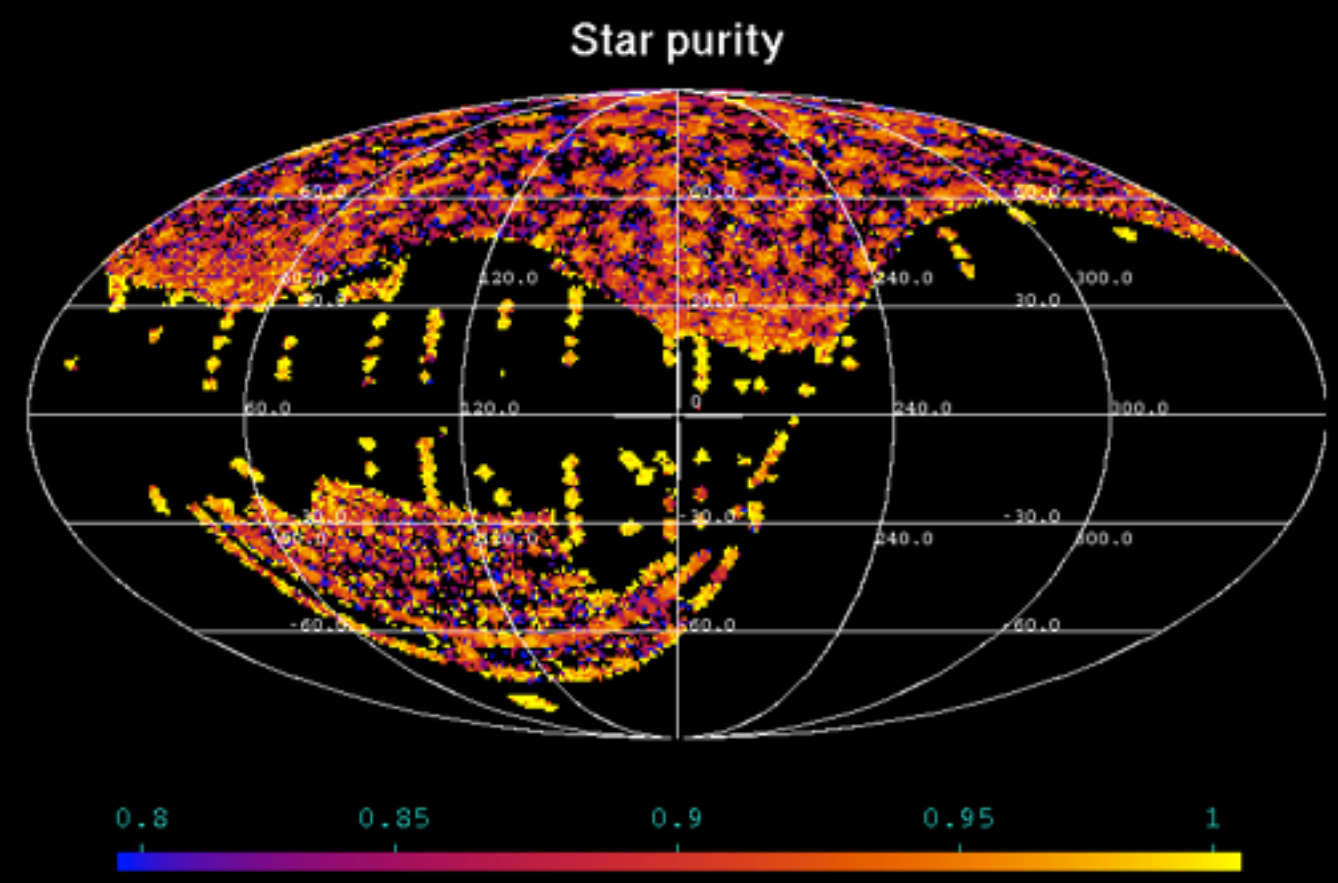
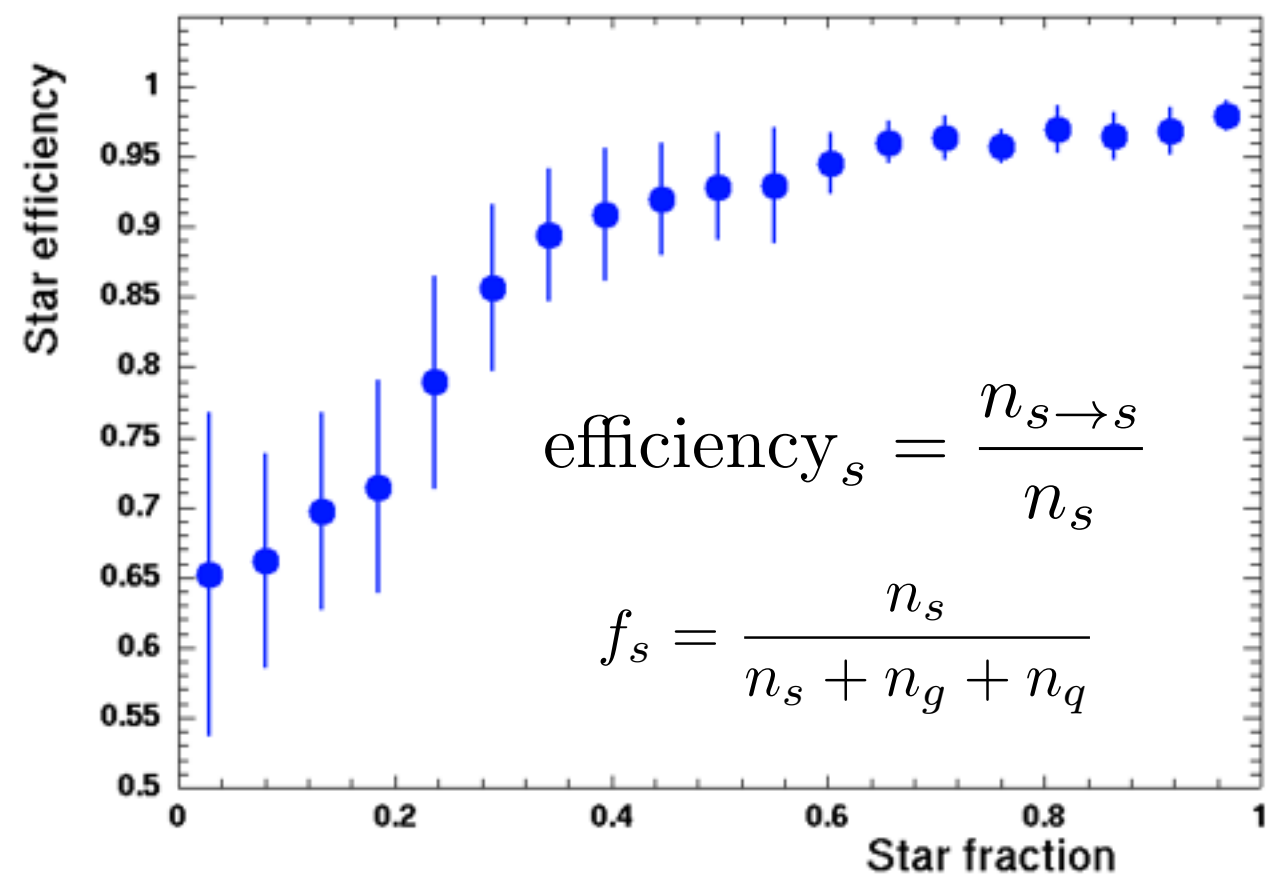
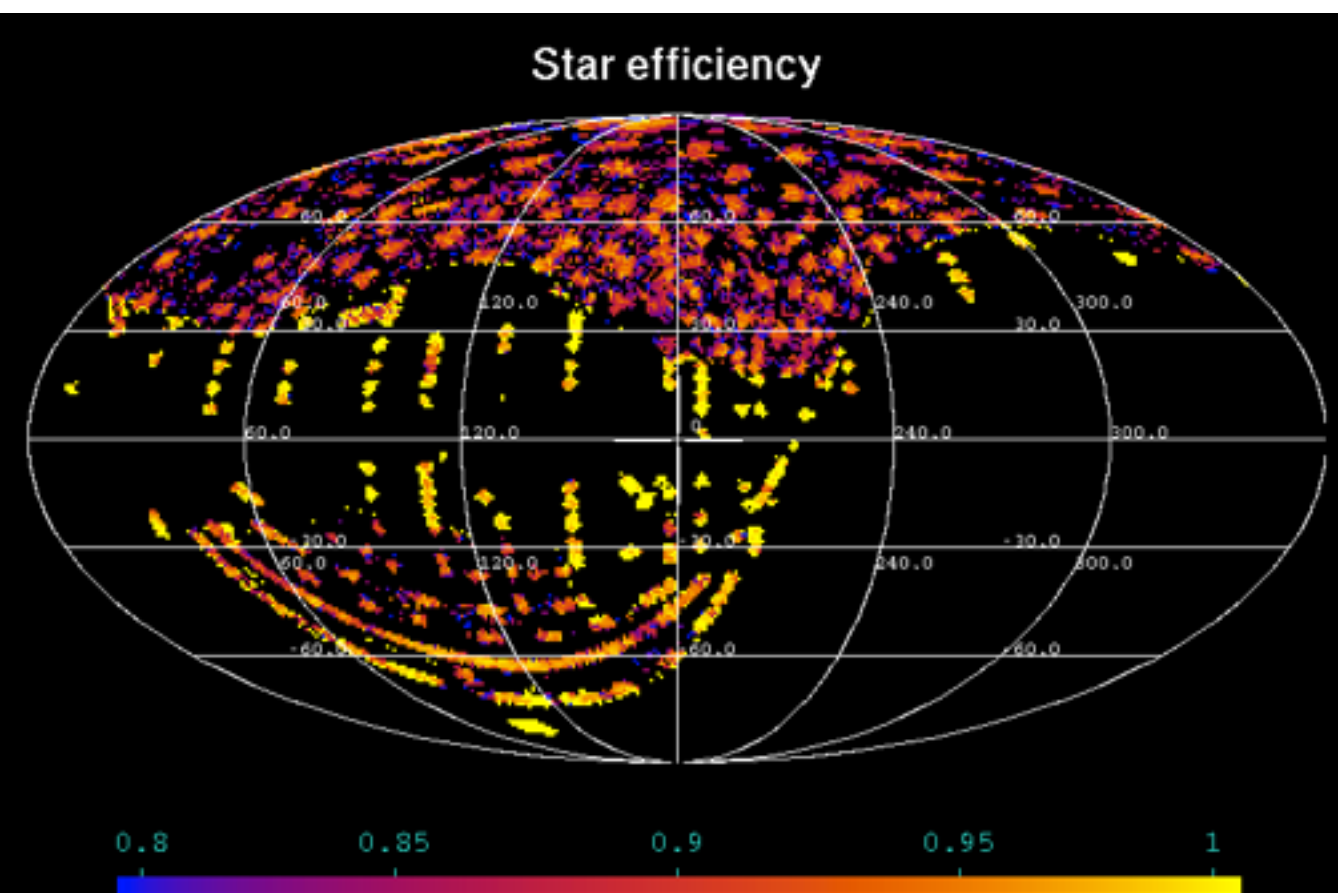


Galaxy efficiency

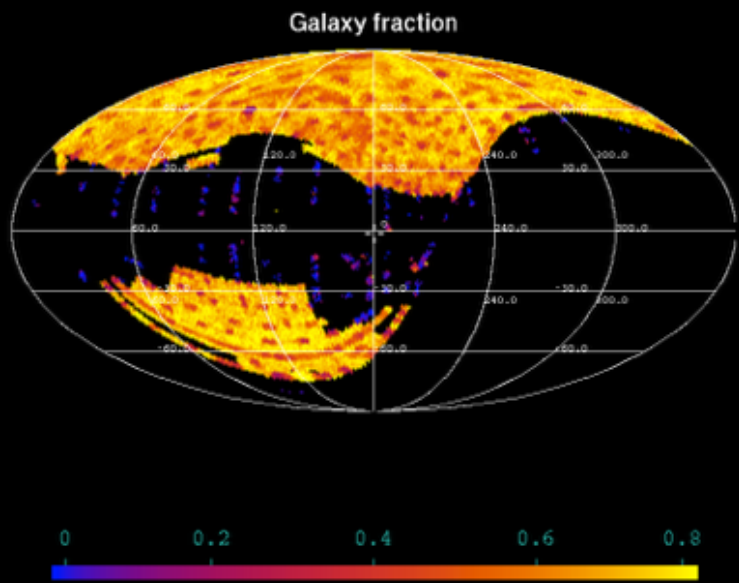


Galaxy purity

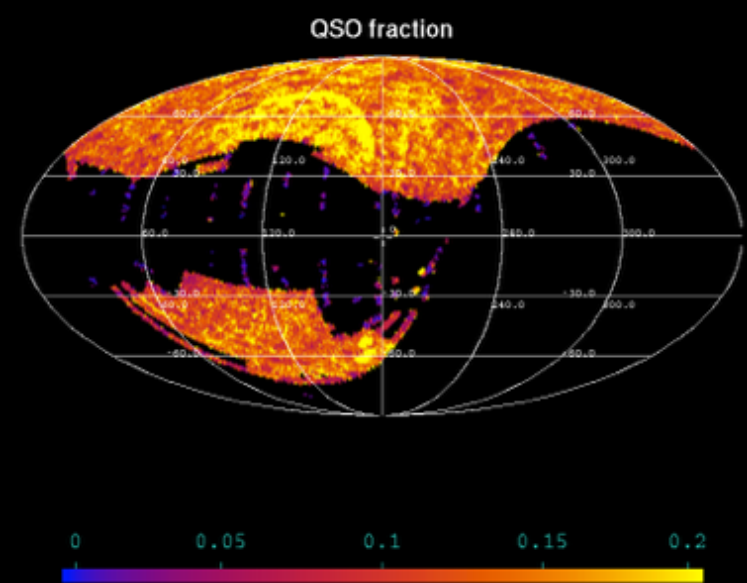
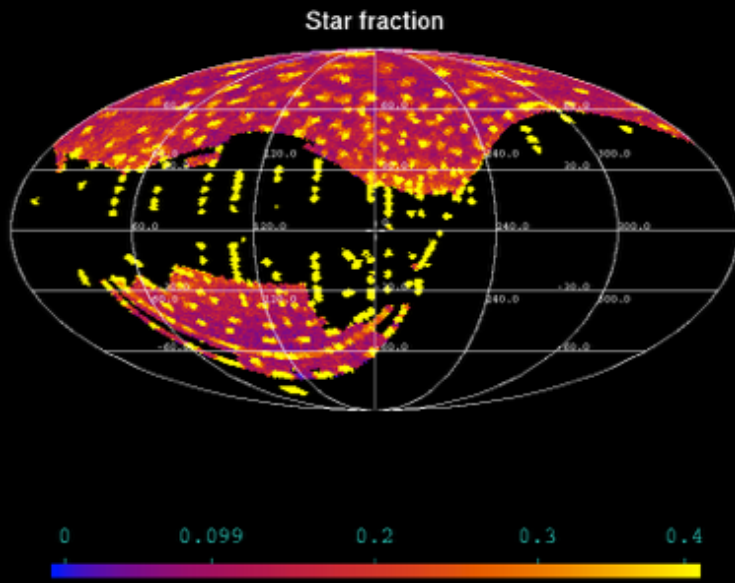






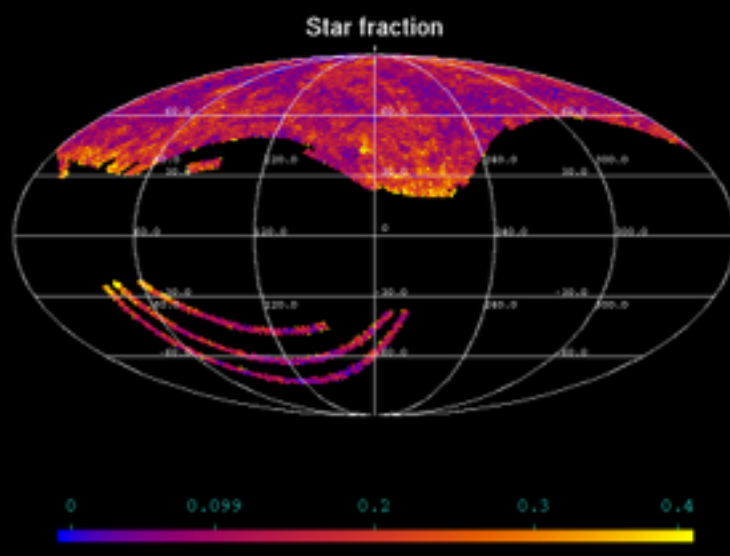
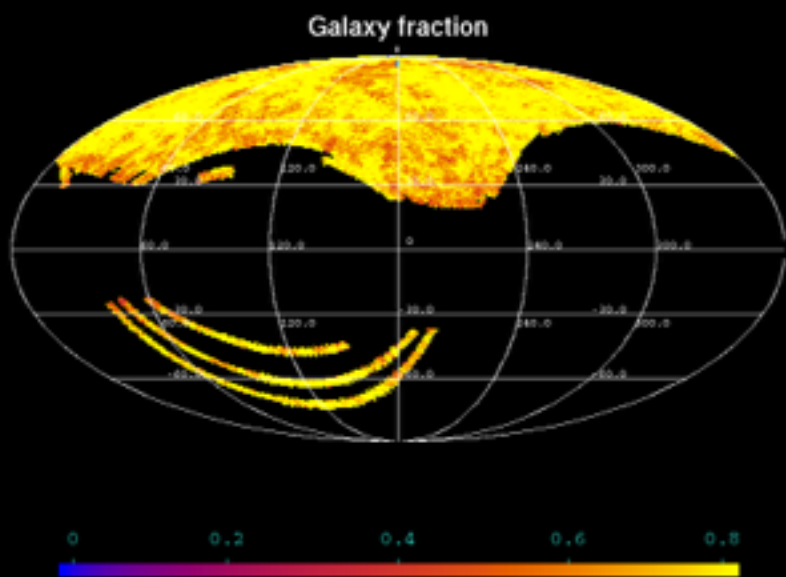


DR 12

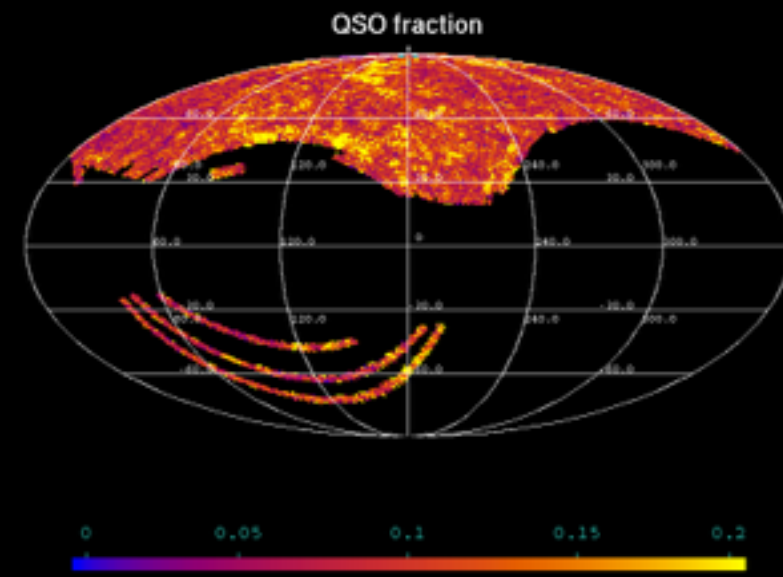


	Stars	Galaxies	QSOs	Total
<b>Num</b>	<b>160,040</b>	<b>879,792</b>	<b>120,425</b>	<b>1,160,257</b>
<b>fraction</b>	<b>14%</b>	<b>76%</b>	<b>10%</b>	
<b>efficiency</b>	<b>95%</b>	<b>99%</b>	<b>90%</b>	<b>98%</b>
<b>purity</b>	<b>94%</b>	<b>99%</b>	<b>94%</b>	

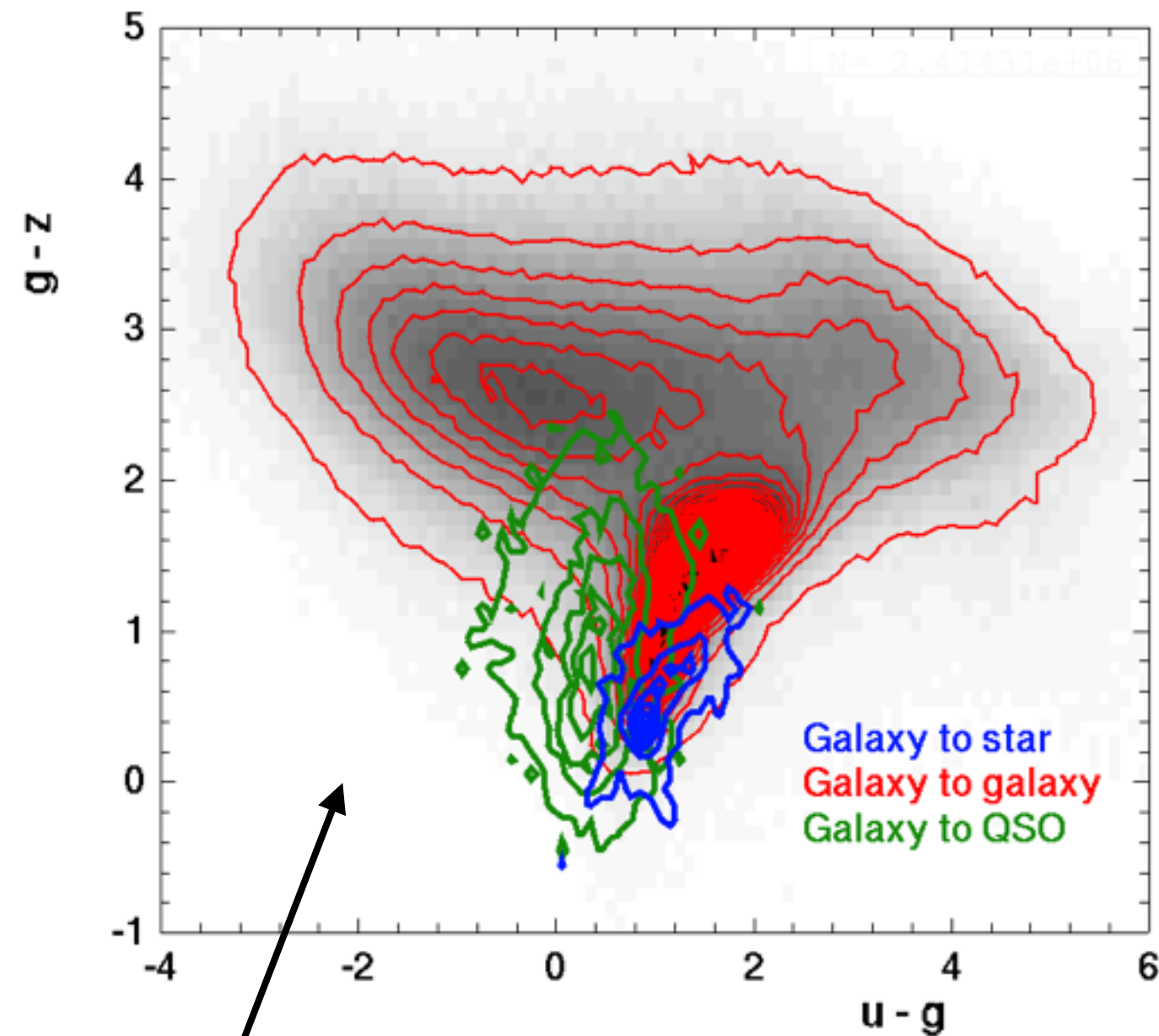
DR12  
vs  
Legacy



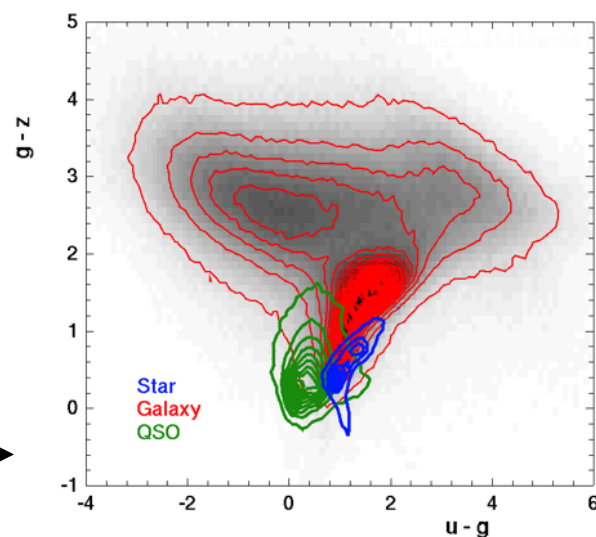
Legacy



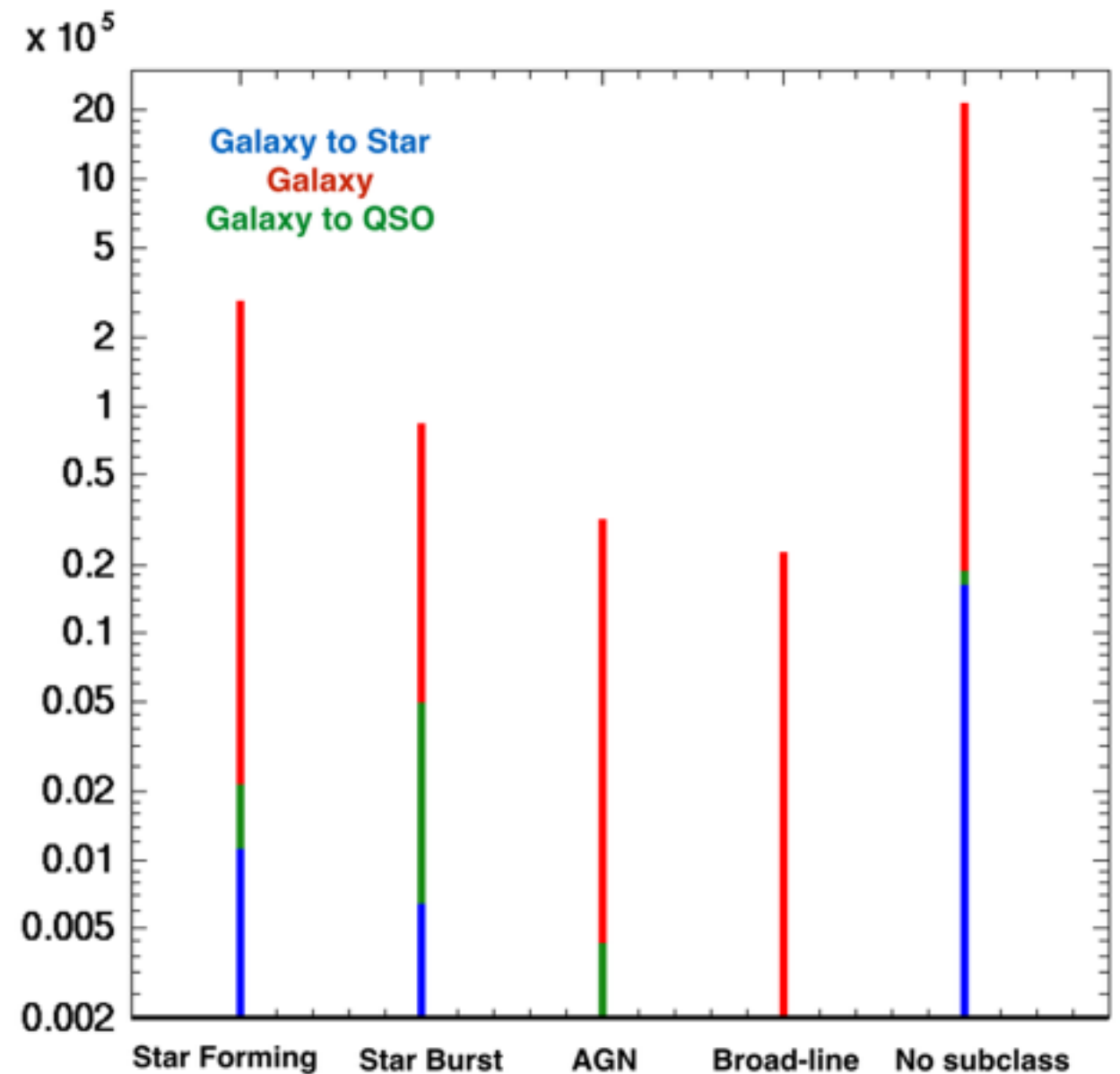
# Galaxy misclassifications



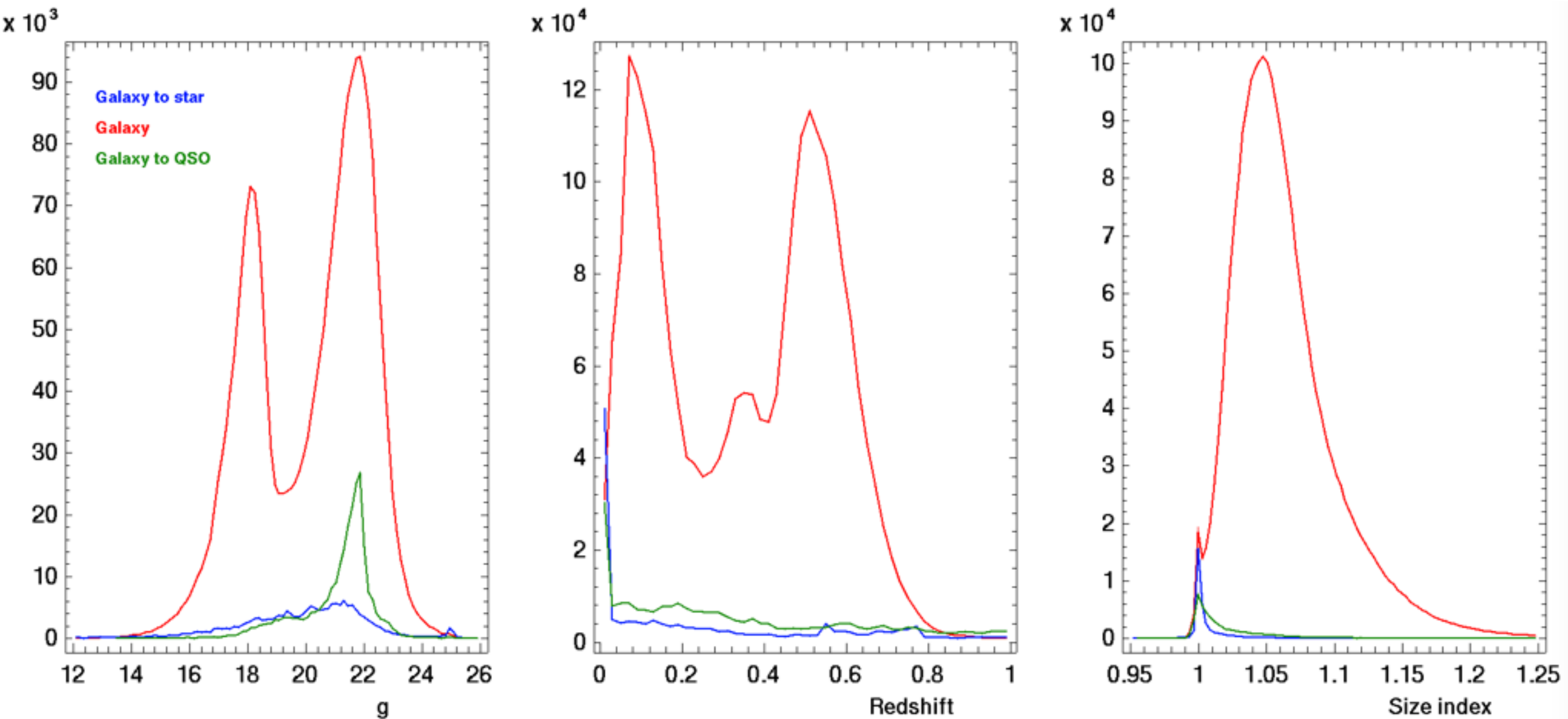
Classified types  
vs  
spectroscopic types



## Misclassified sub-classes



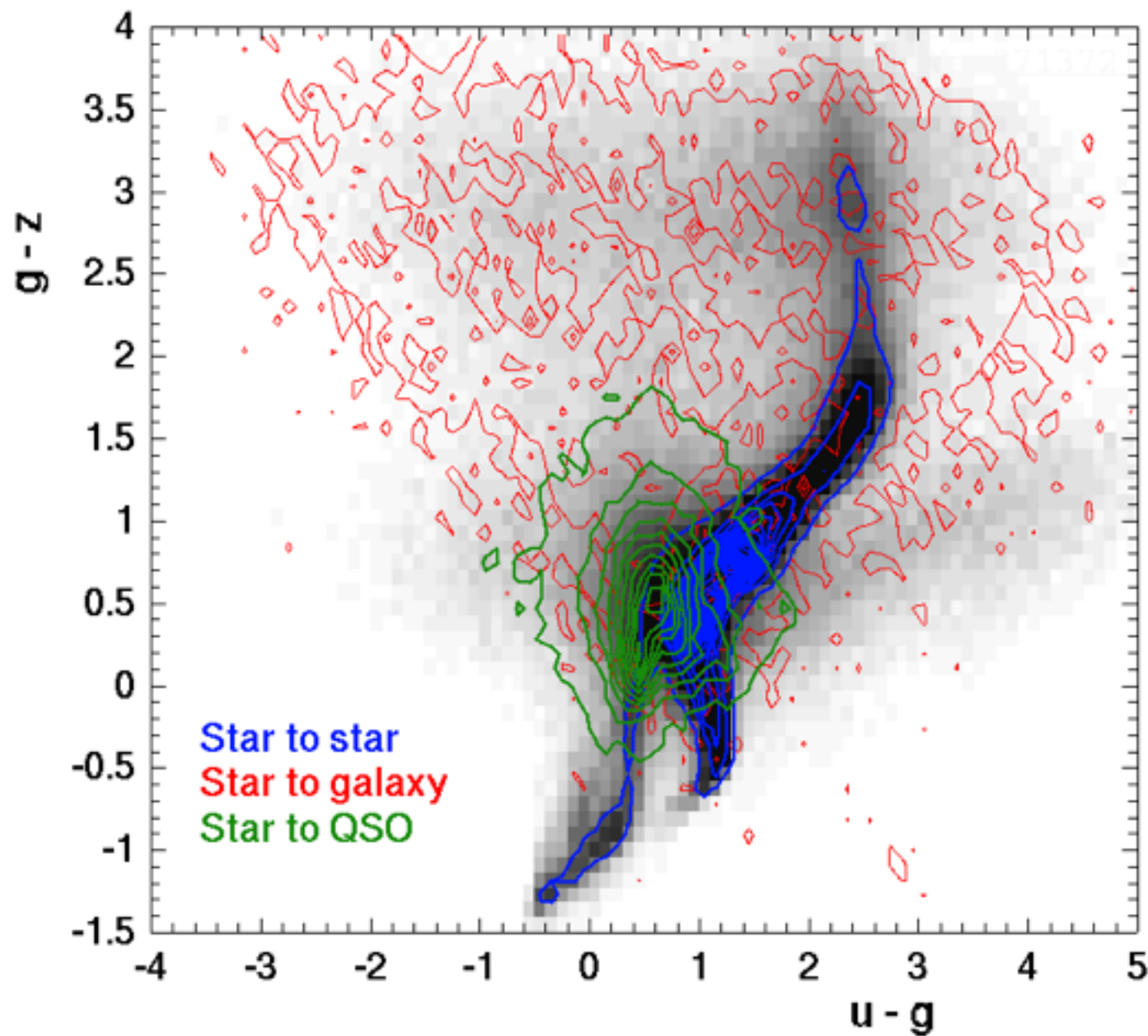
# Galaxy misclassifications



Faint nearby galaxies ==> point-like sources

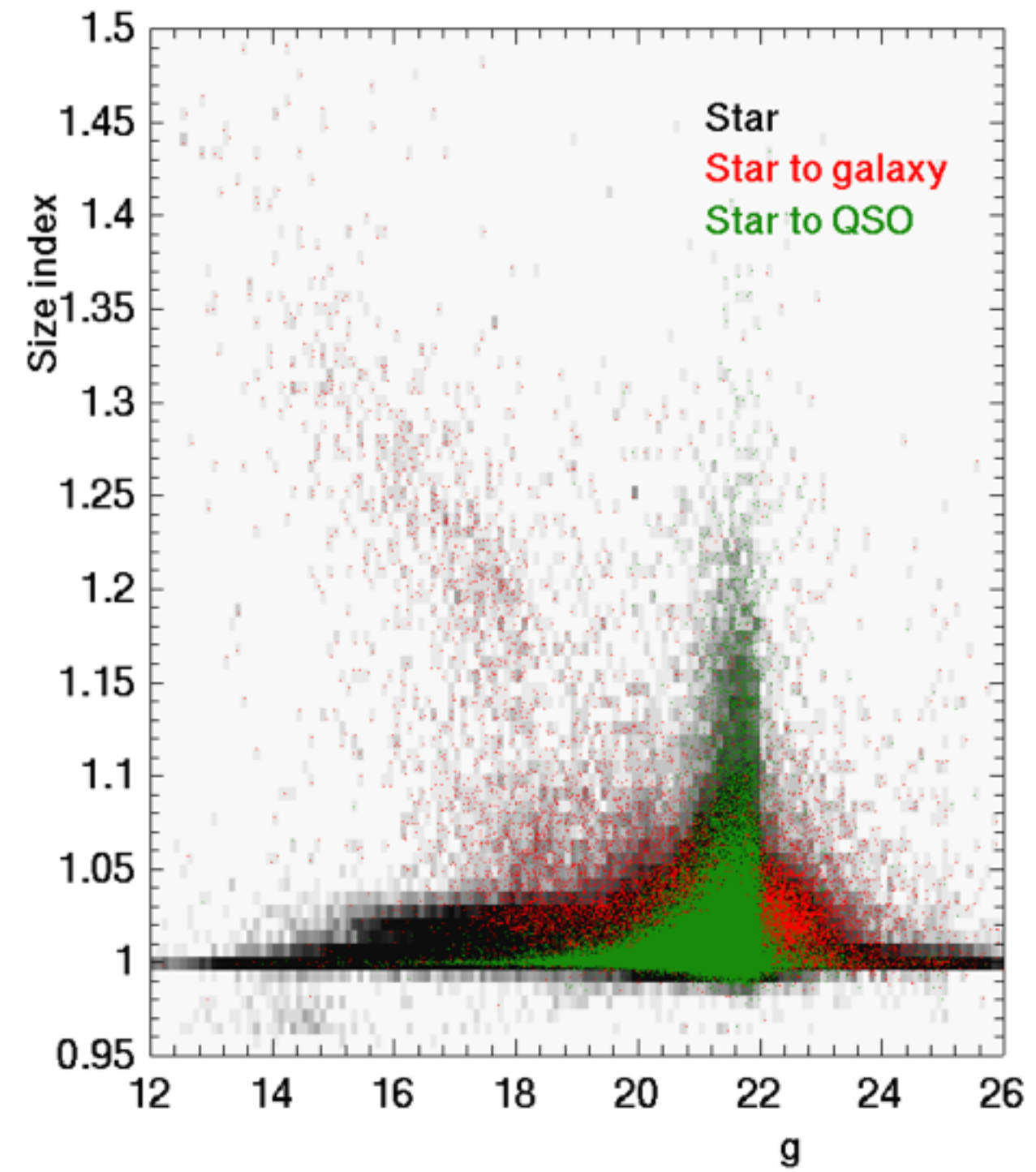


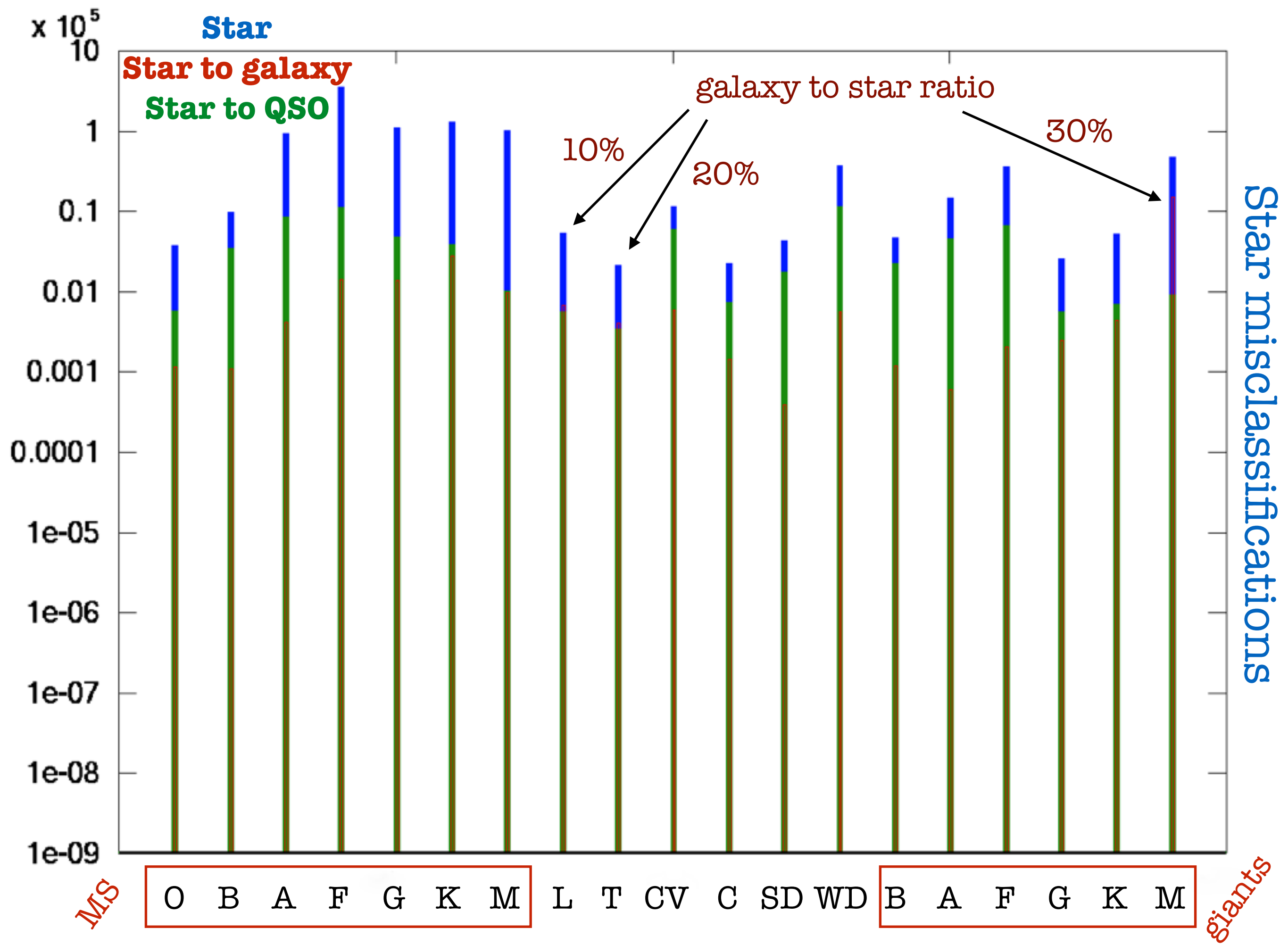
# Star misclassifications



stars with  
scattered colours  
contaminates galaxy sample

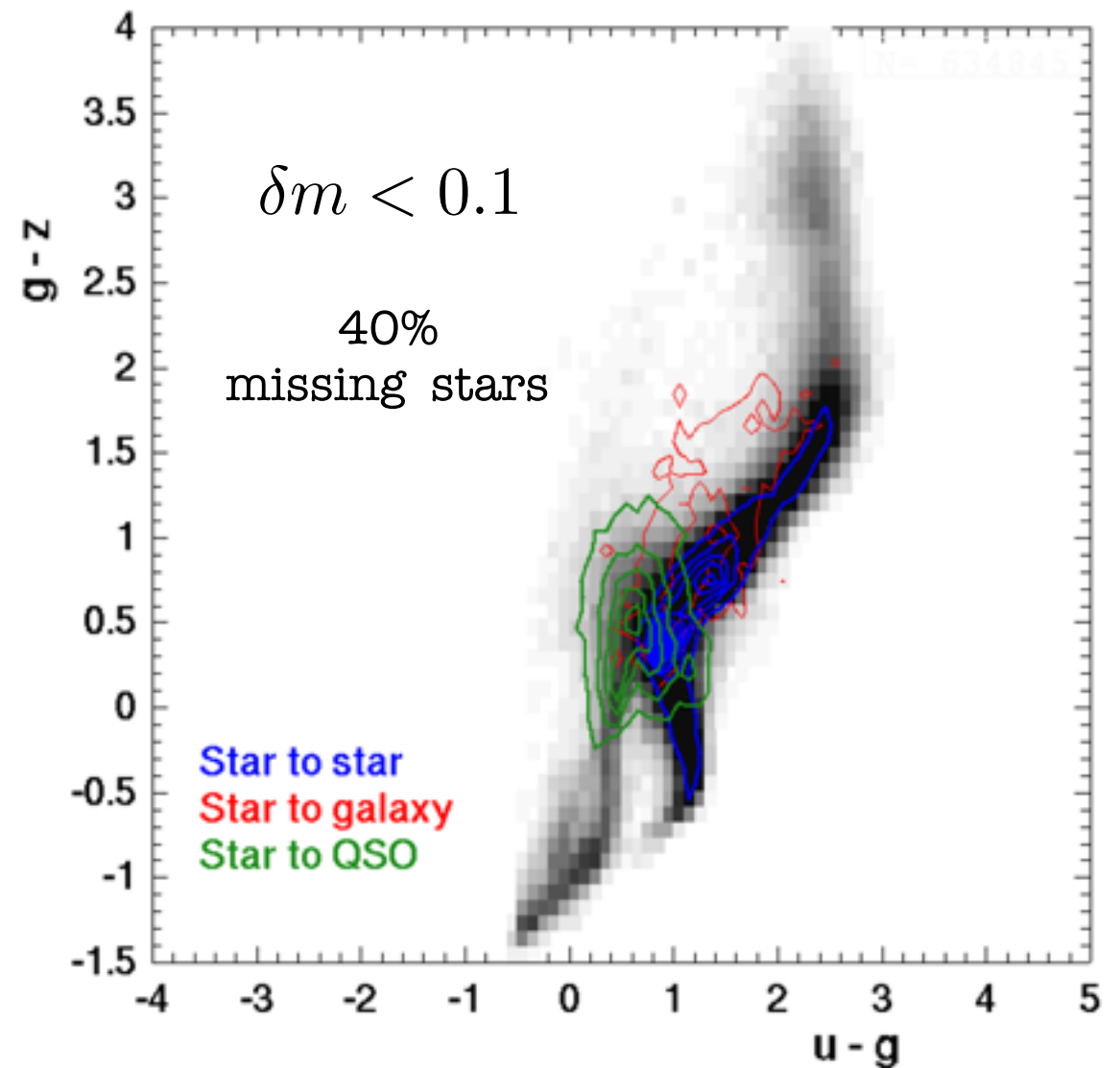
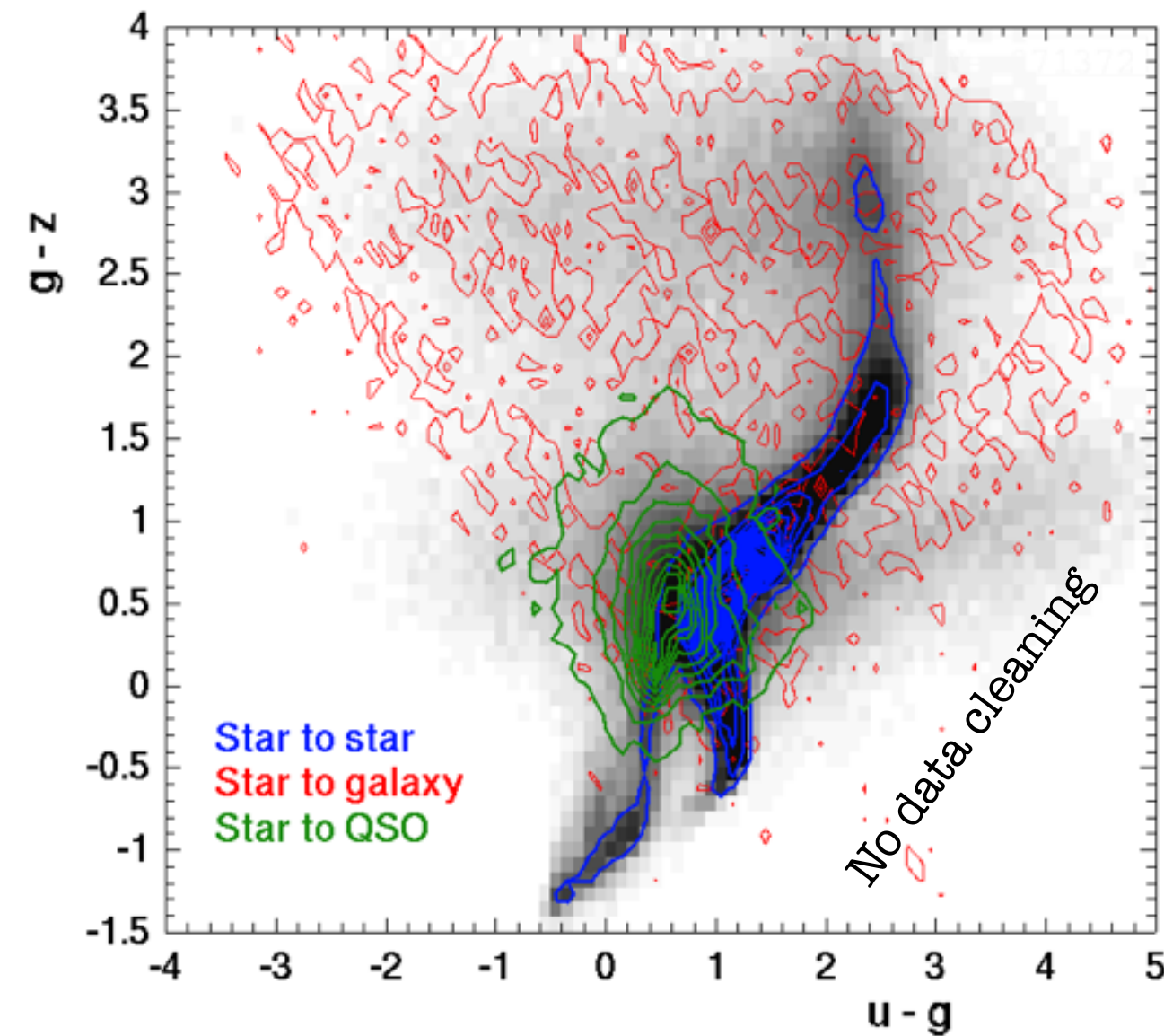
deviation from point-like source  
for faint stars





# Star misclassifications

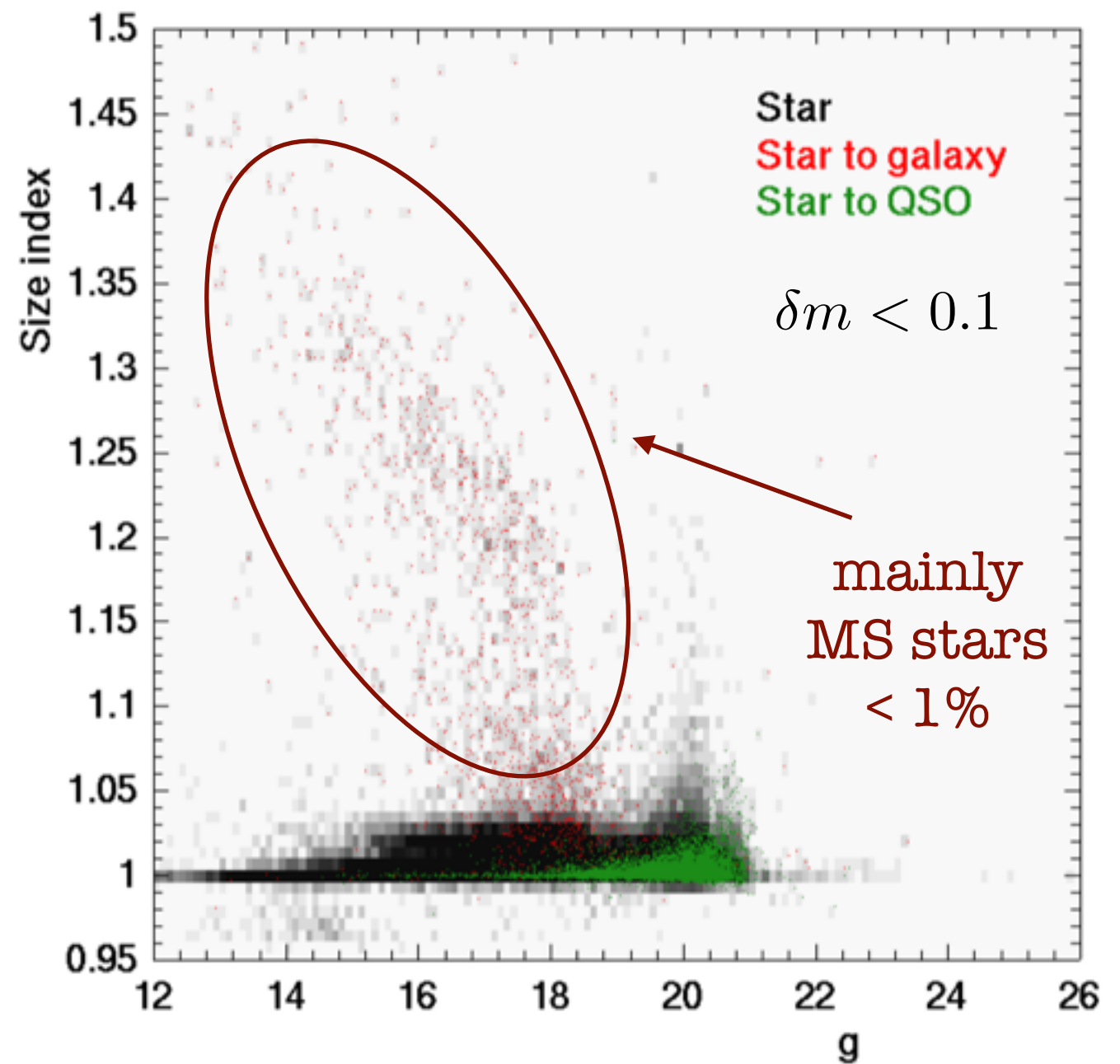
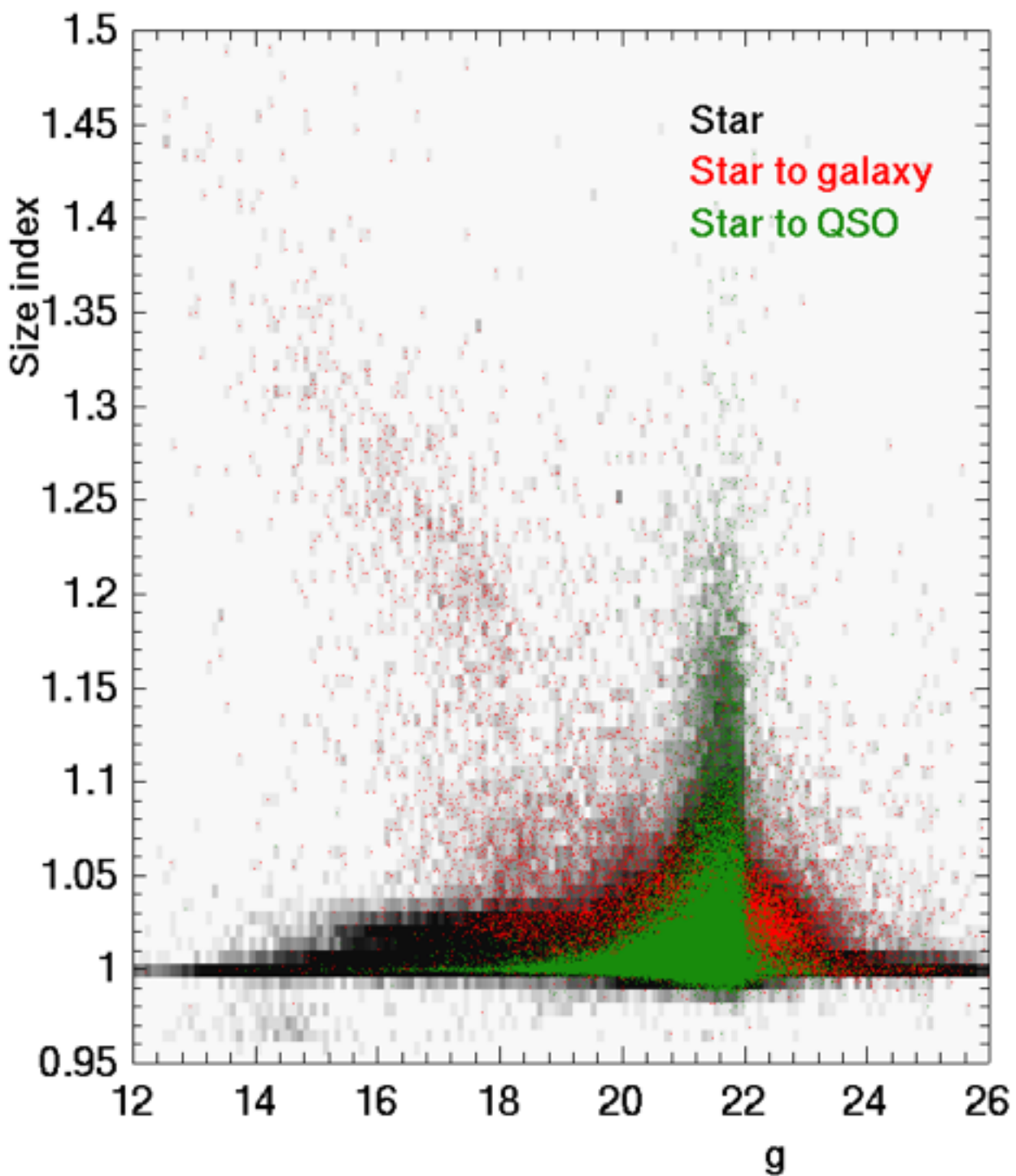
Photometric quality affects the colour measurement accuracy



	Star	Galaxy	QSO	Total
Efficiency	97.1%	97.6%	93.8%	97%
Purity	97.5%	98.9%	94.9%	-



# Star misclassifications



# Summary

- Colour indices and apparent angular size can separates galaxies from stars and QSOs
- For SDSS DR12 a 4-layer MLP separates galaxies from the point-like sources by precision better than 98%
- Observational strategy with uniform sky coverage improves the classification efficiency
- Faint nearby galaxies can be misclassified as point-like sources while redshifted galaxies tend less to be misclassified
- M-giant stars, faint red L and T stars mainly contaminate the classified galaxy sample

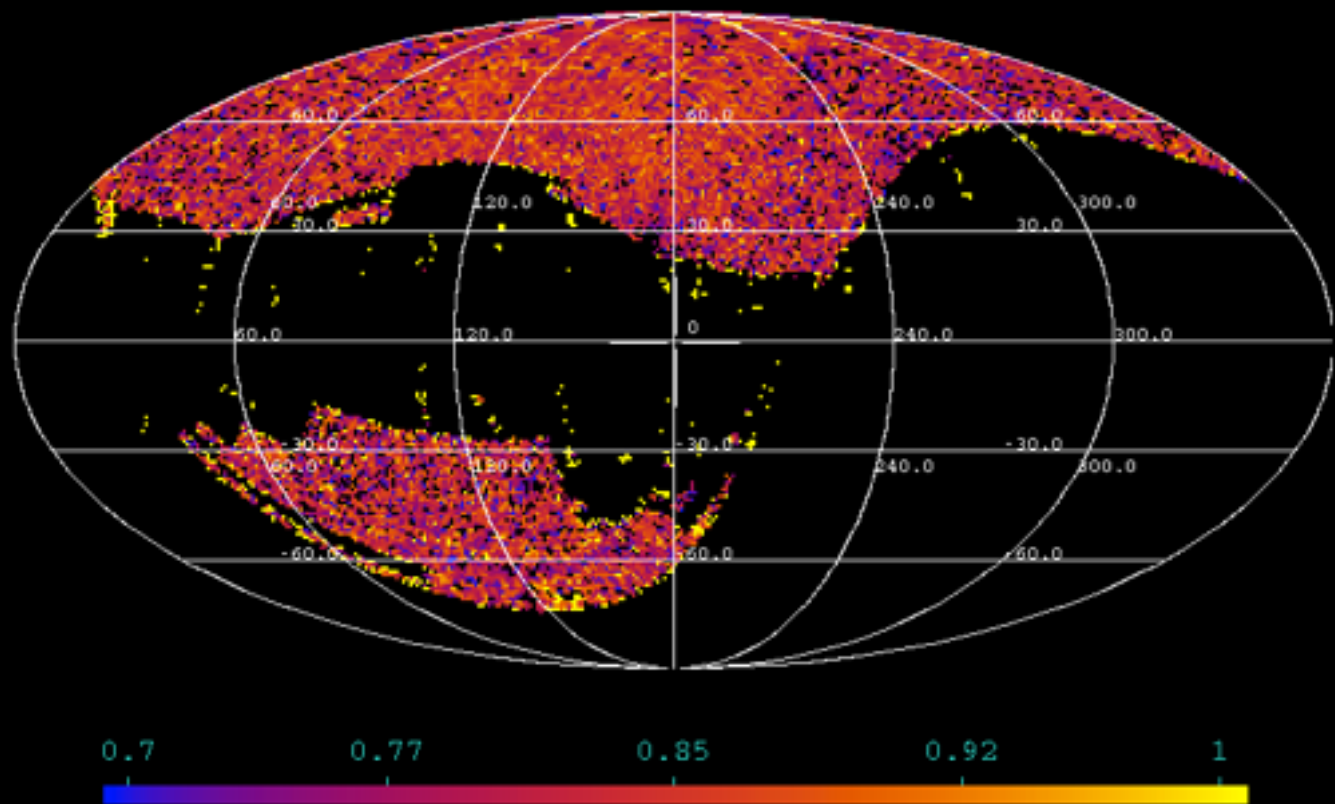
# Classification for the LSST

- Generating different galaxy types according to their luminosity function and the LSST apparent magnitude limits
- Simulating the colour indices according to galaxy redshifted SEDs and LSST pass-band filters
- Including the stars

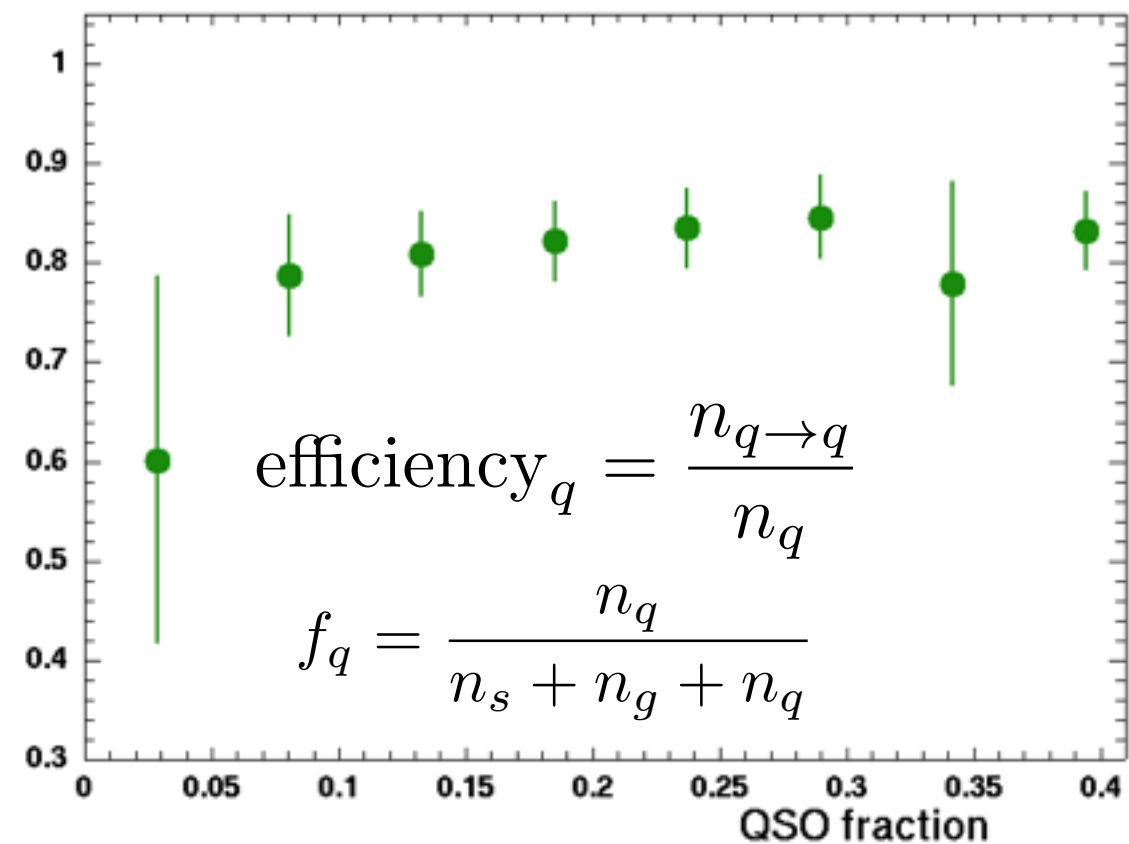


# Backups

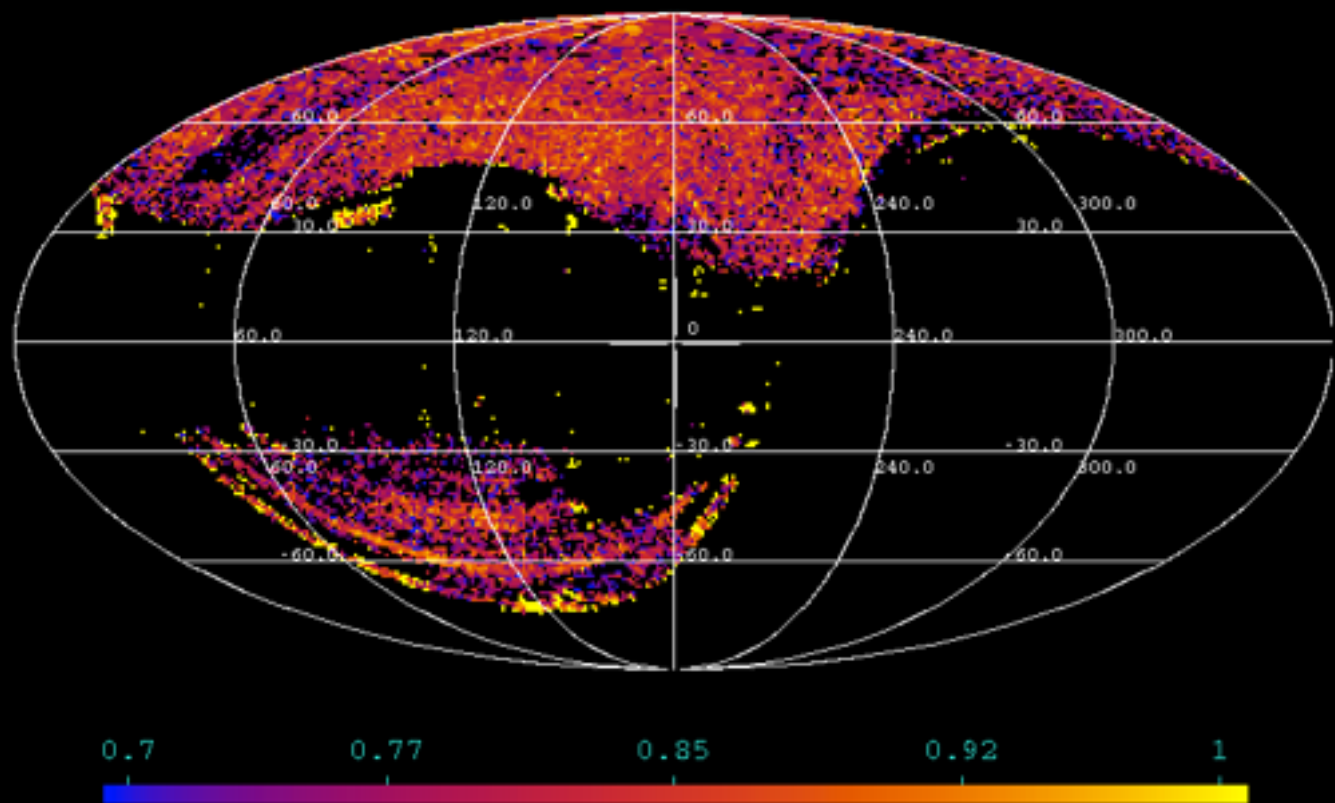
QSO efficiency



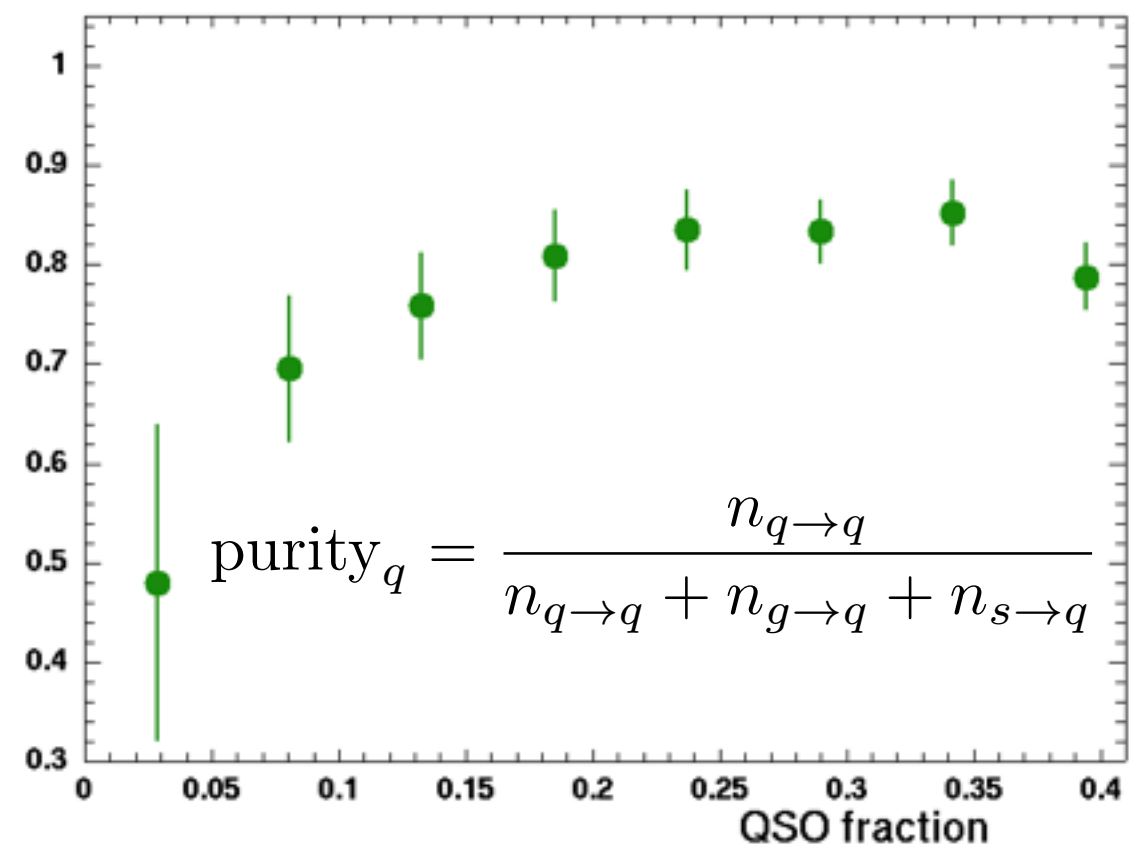
QSO efficiency



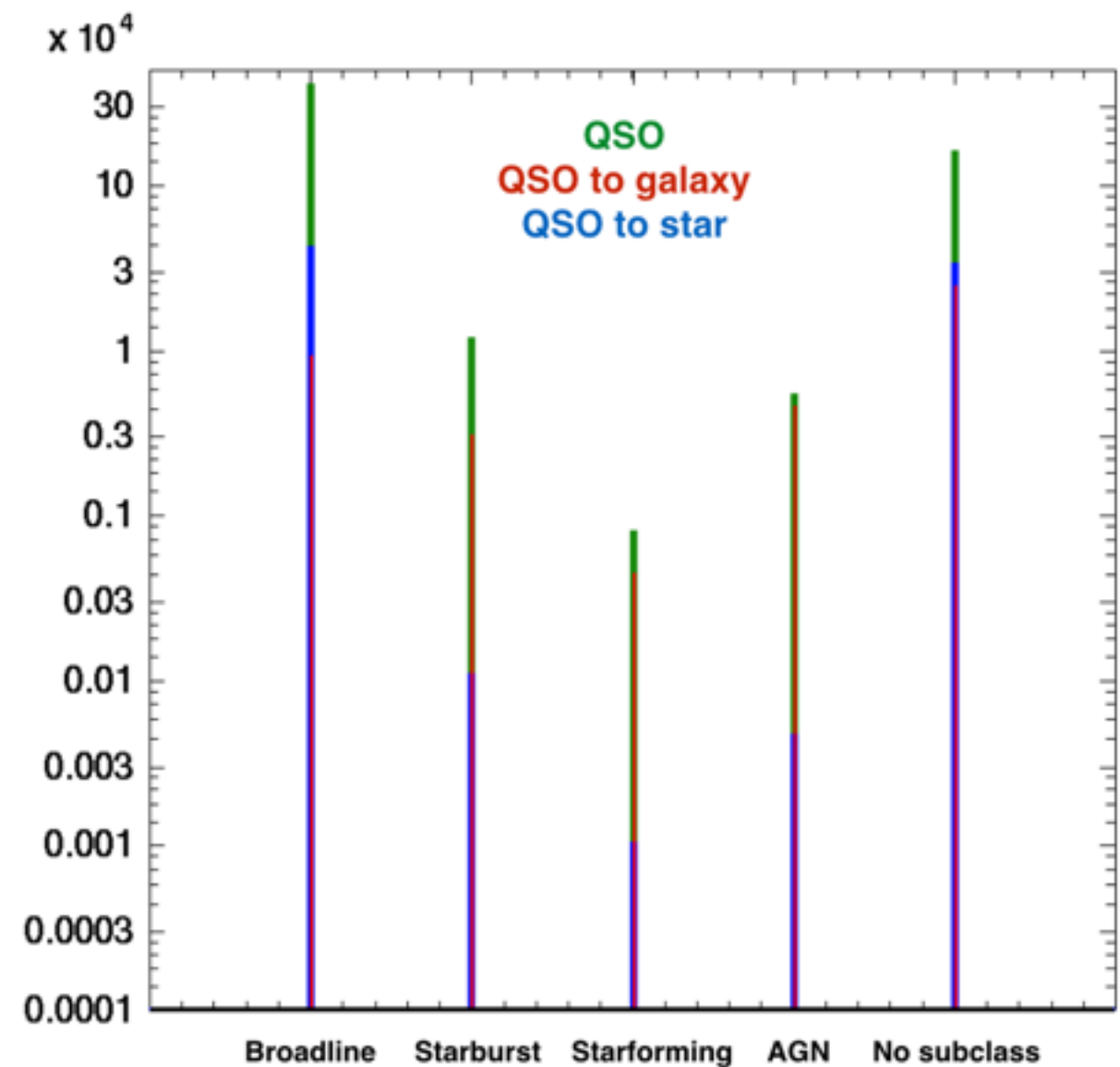
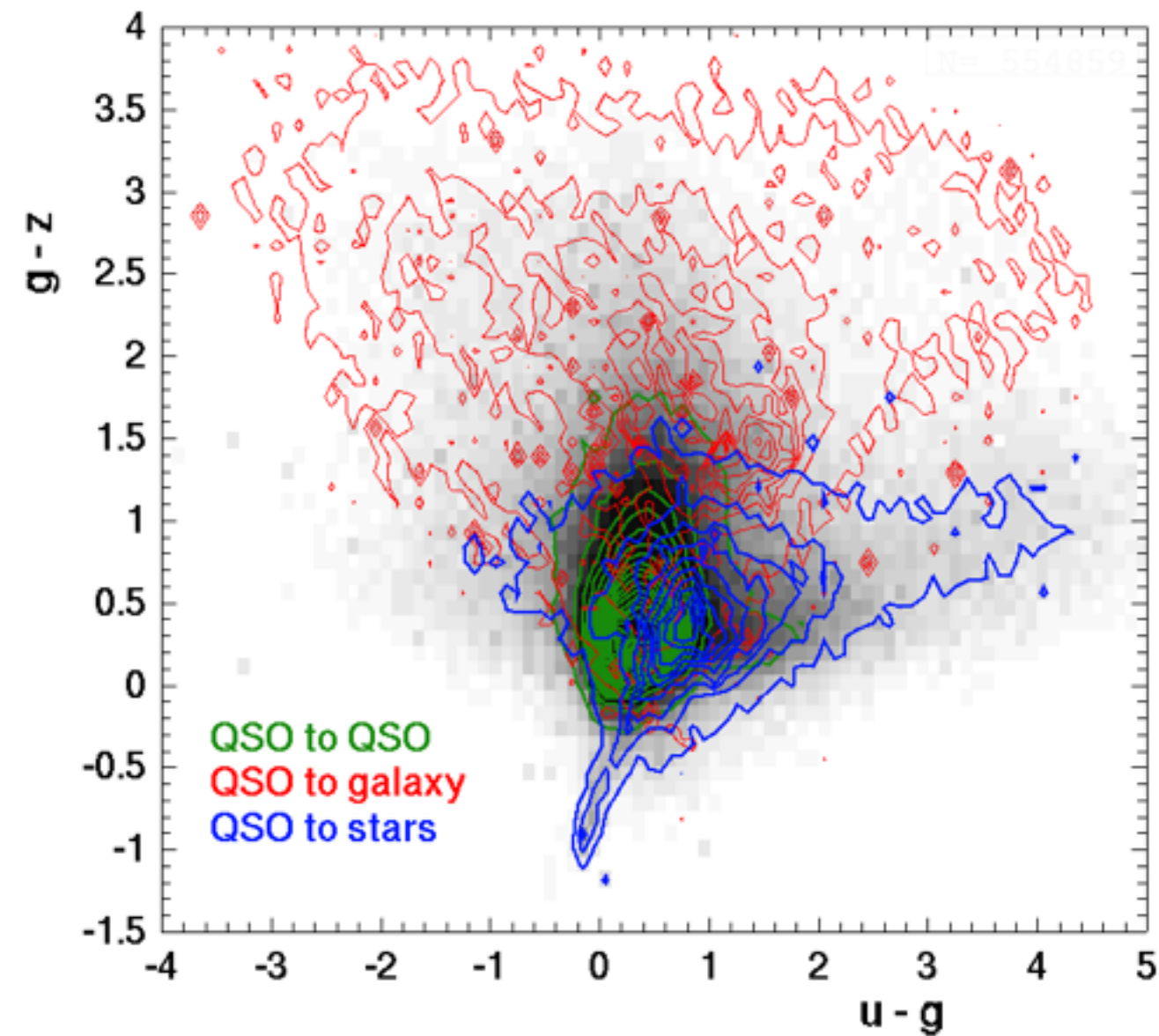
QSO purity



QSO purity

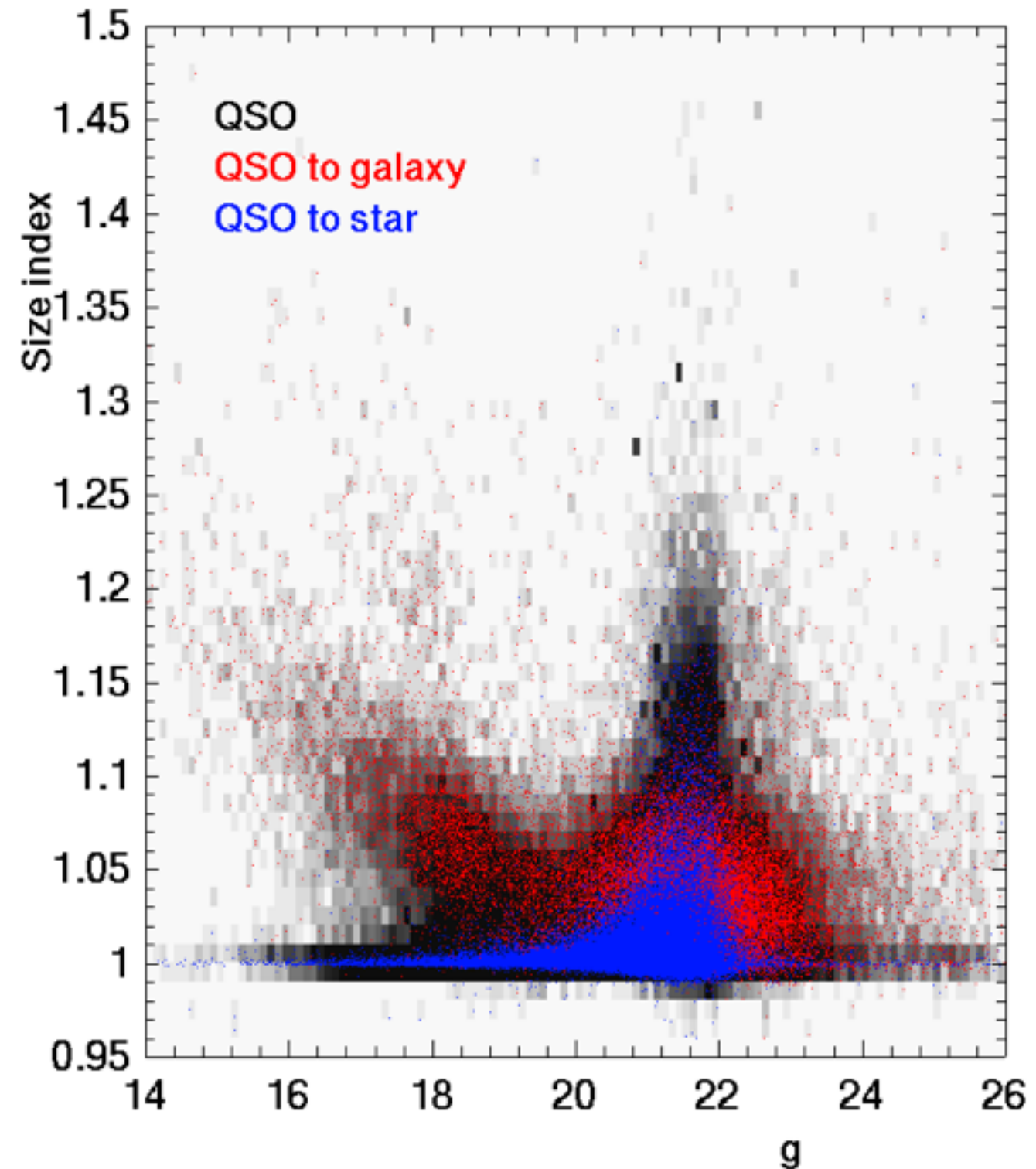
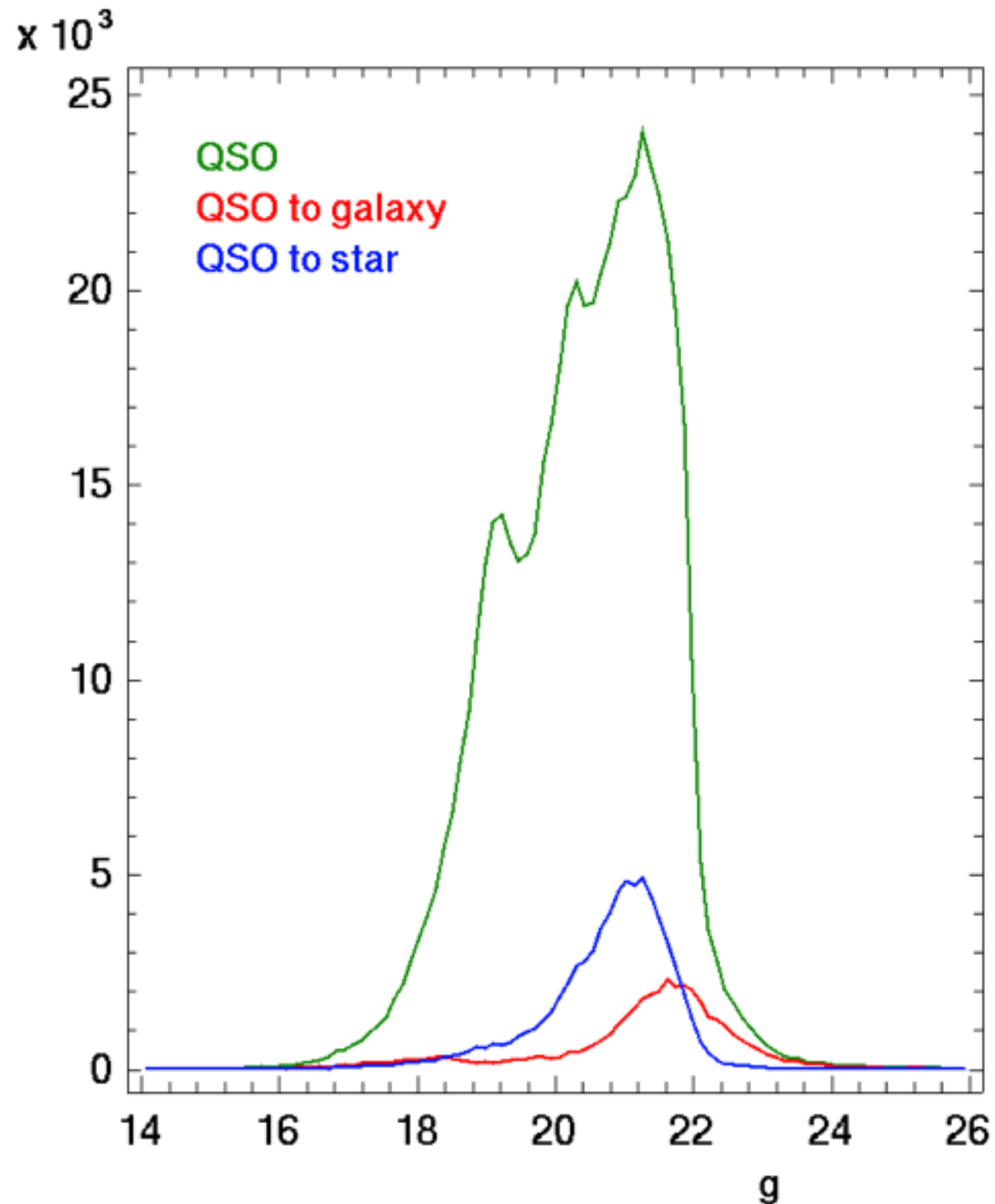


# QSOs misclassifications

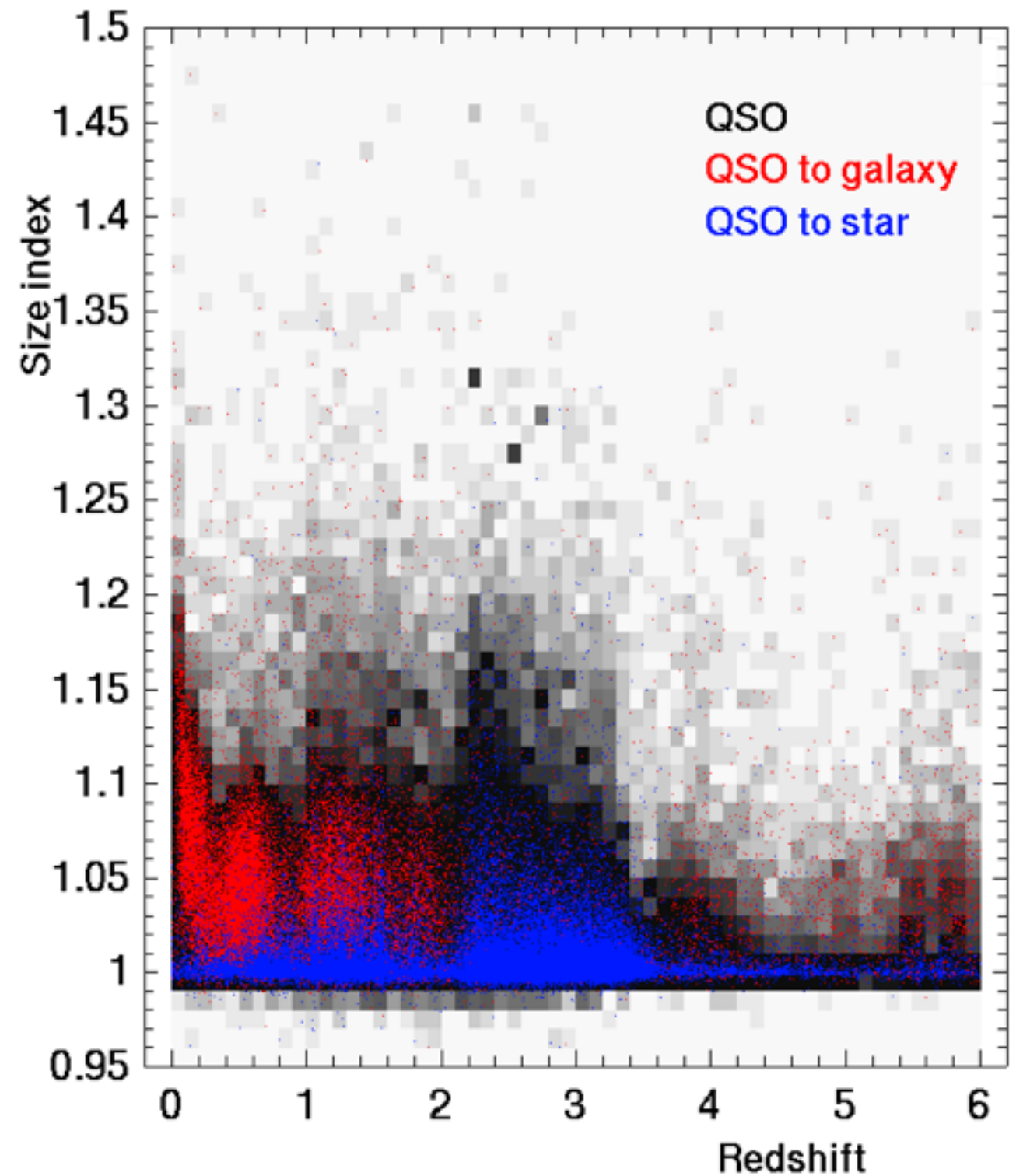
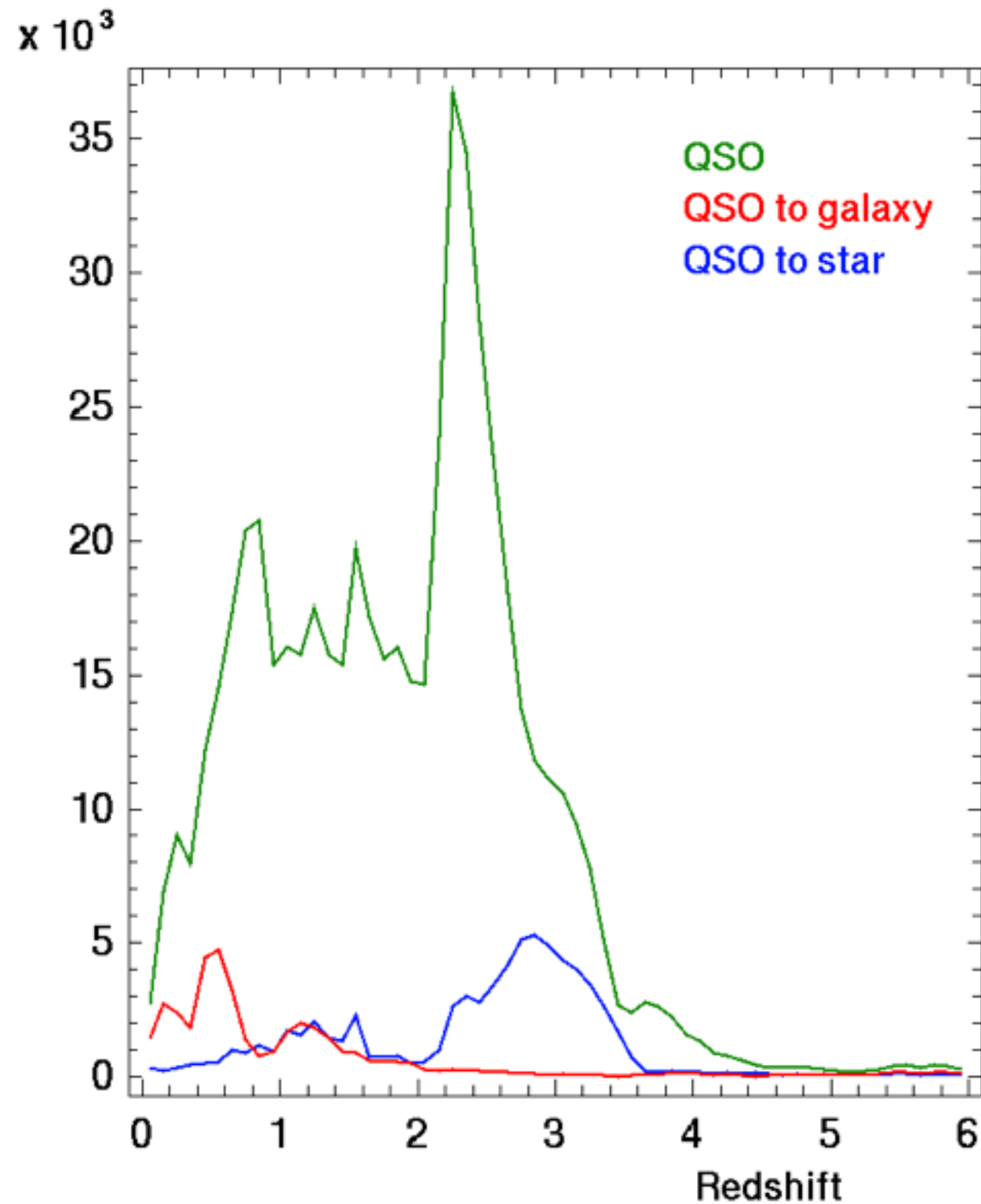


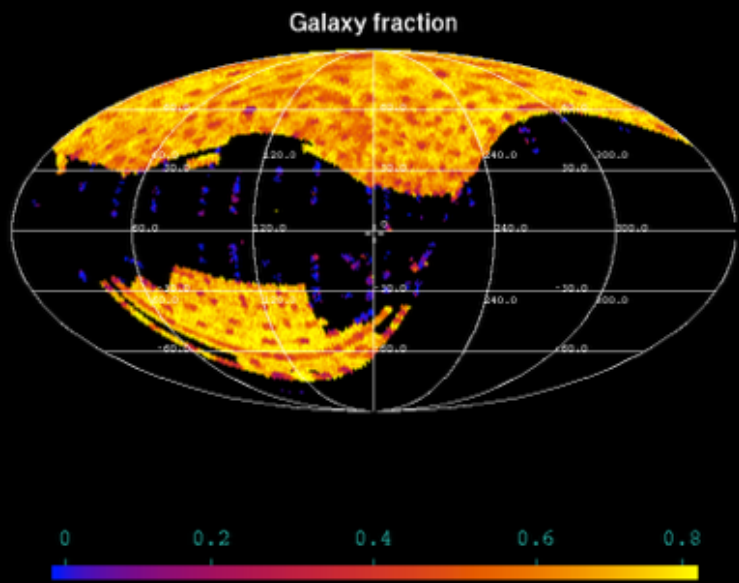


# QSOs misclassifications

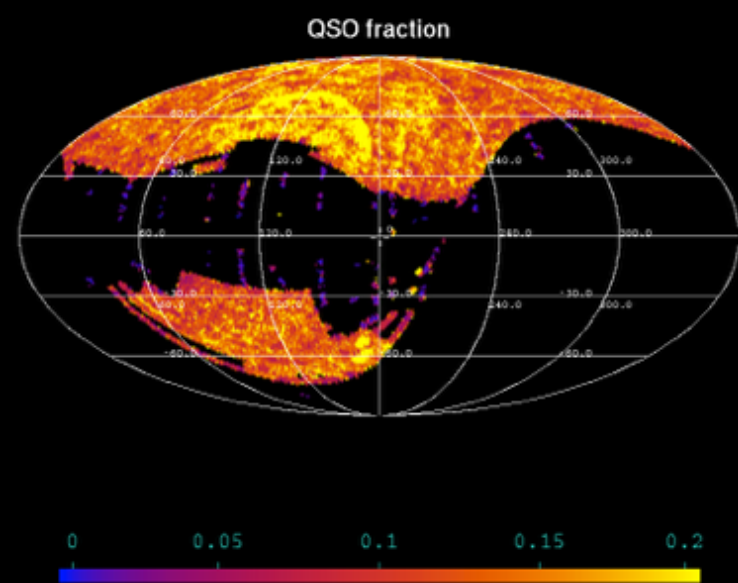
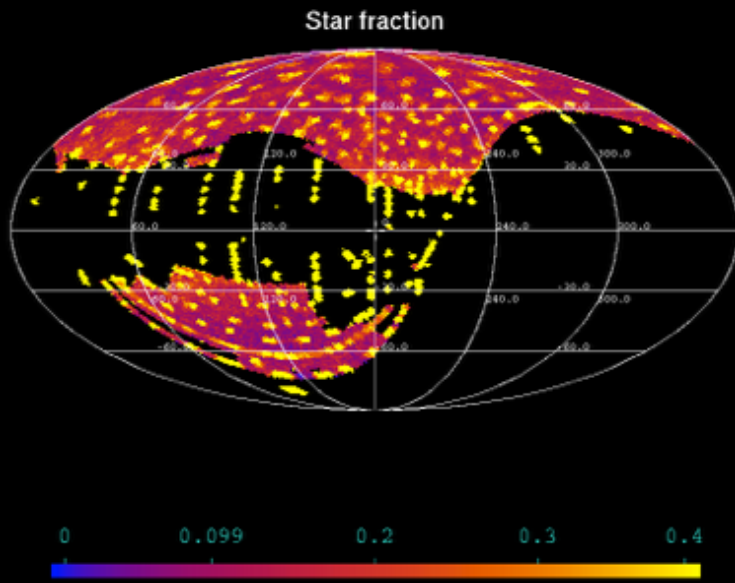


# QSOs misclassifications



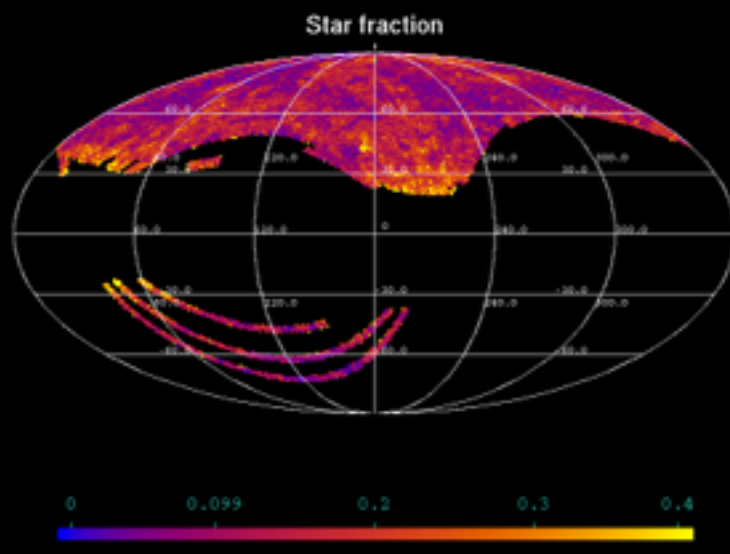
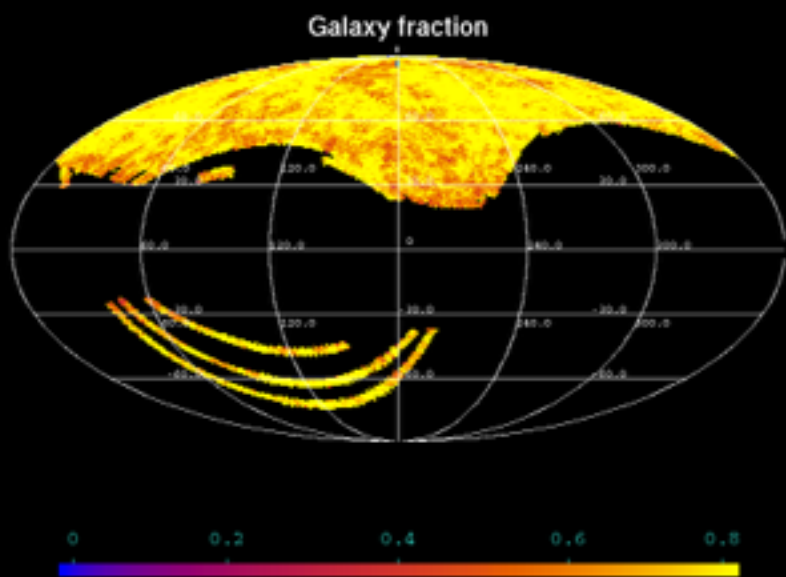


DR 12

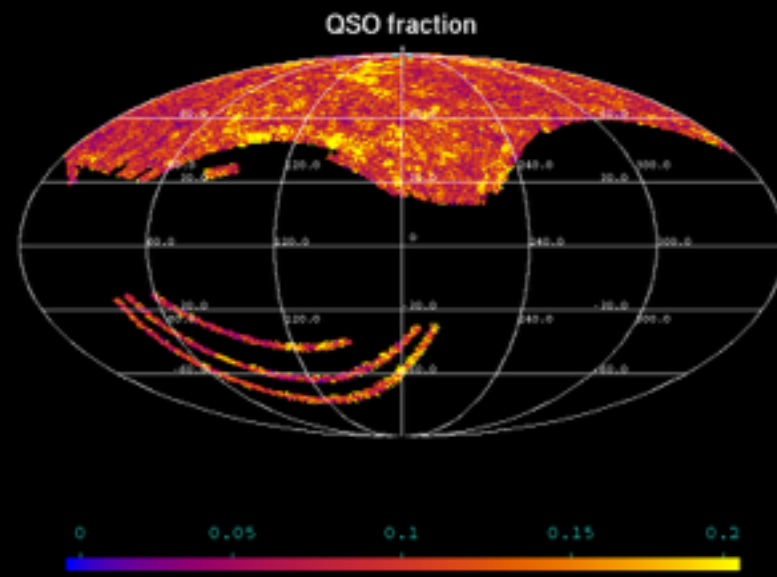


	Stars	Galaxies	QSOs	Total
efficiency	95%	99%	90%	98%
purity	94%	99%	94%	
efficiency	90%	99%	90%	97%
purity	96%	98%	96%	

DR12  
vs  
Legacy

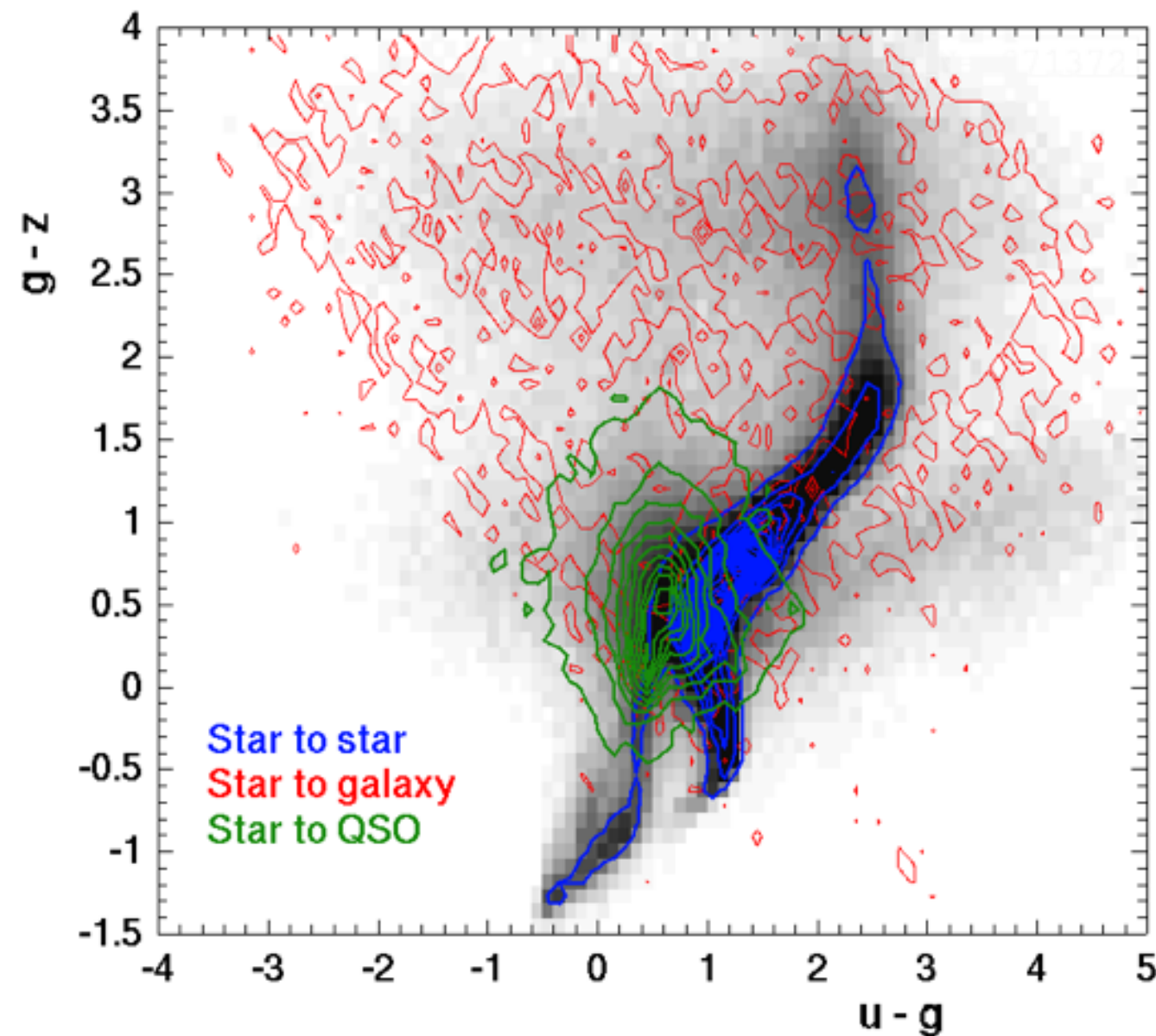


Legacy





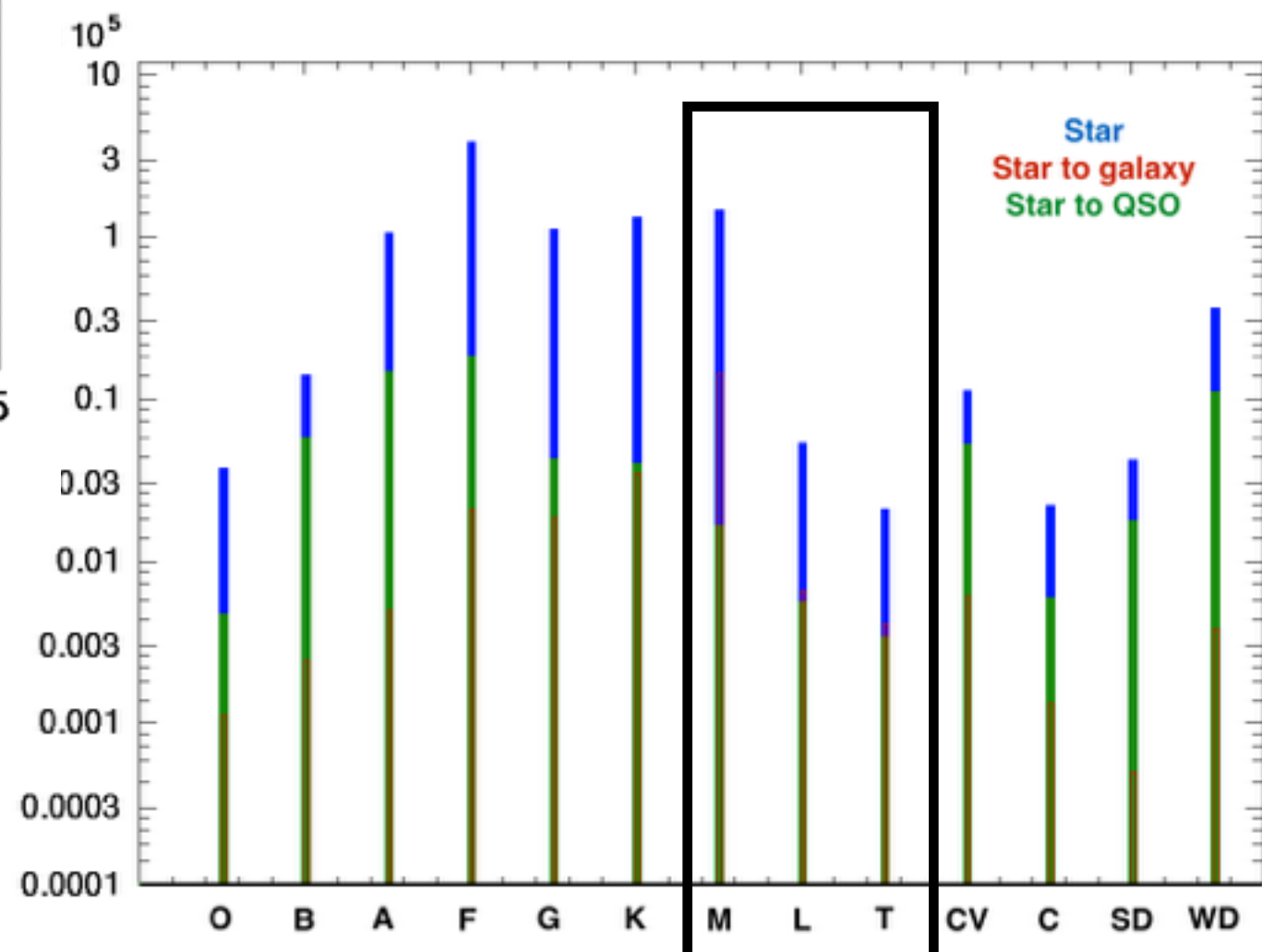
# Star misclassifications



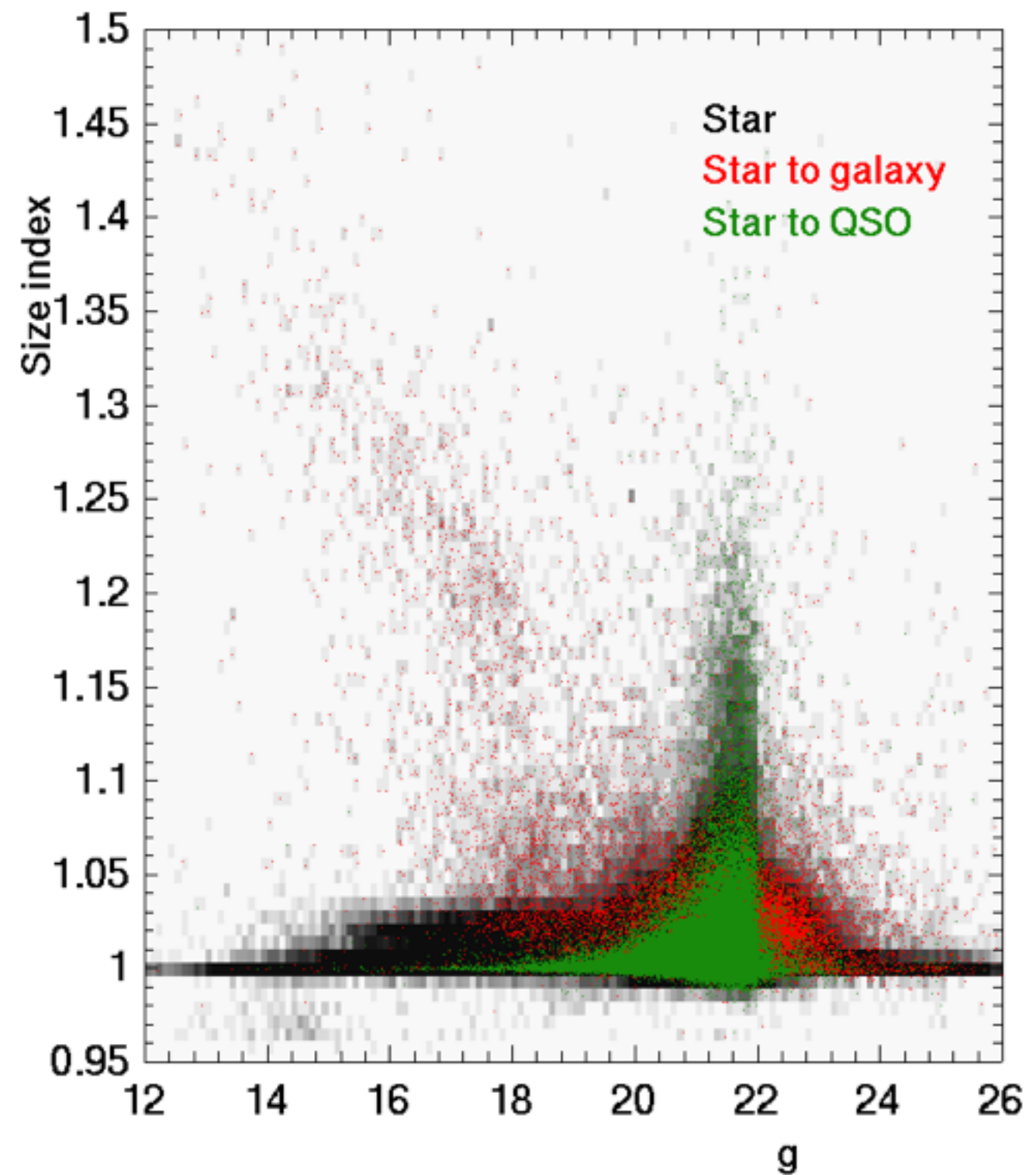
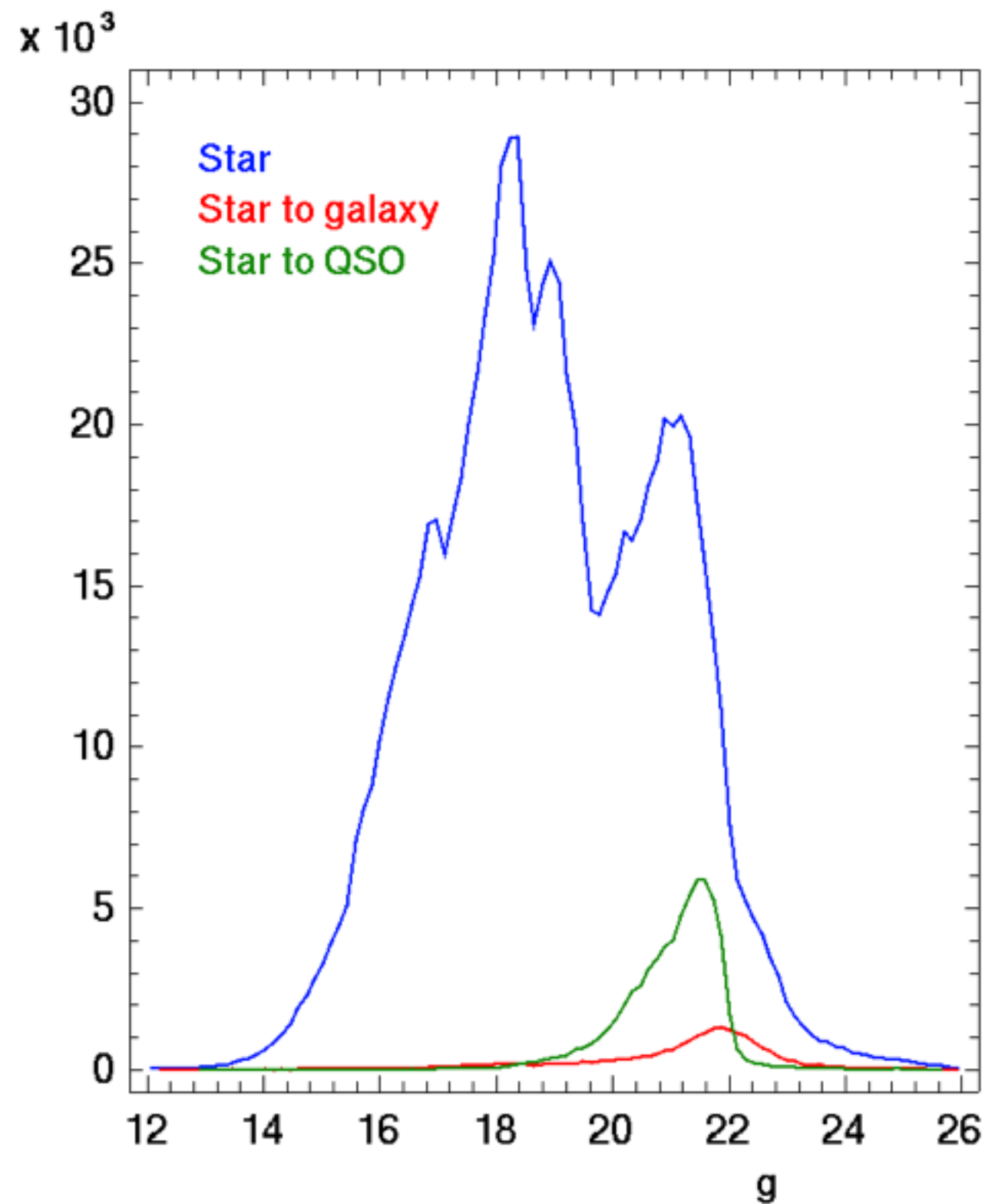
stars with  
scattered colours  
contaminates galaxy sample

Misclassified sub-classes

> 10% of M, L, T stars ==> galaxy



# Star misclassifications



deviation from point-like source for faint stars