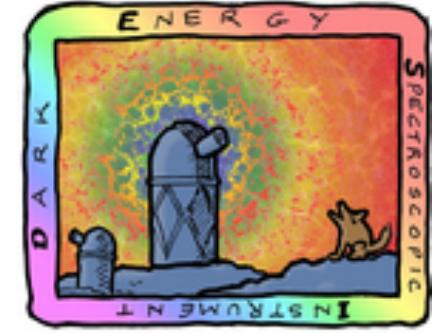


Target Selection of QSO in DESI with Random Forest



E. Burtin, C.-A. Claveau,
N. Palanque-Delabrouille, Ch. Yèche

AI telecon,

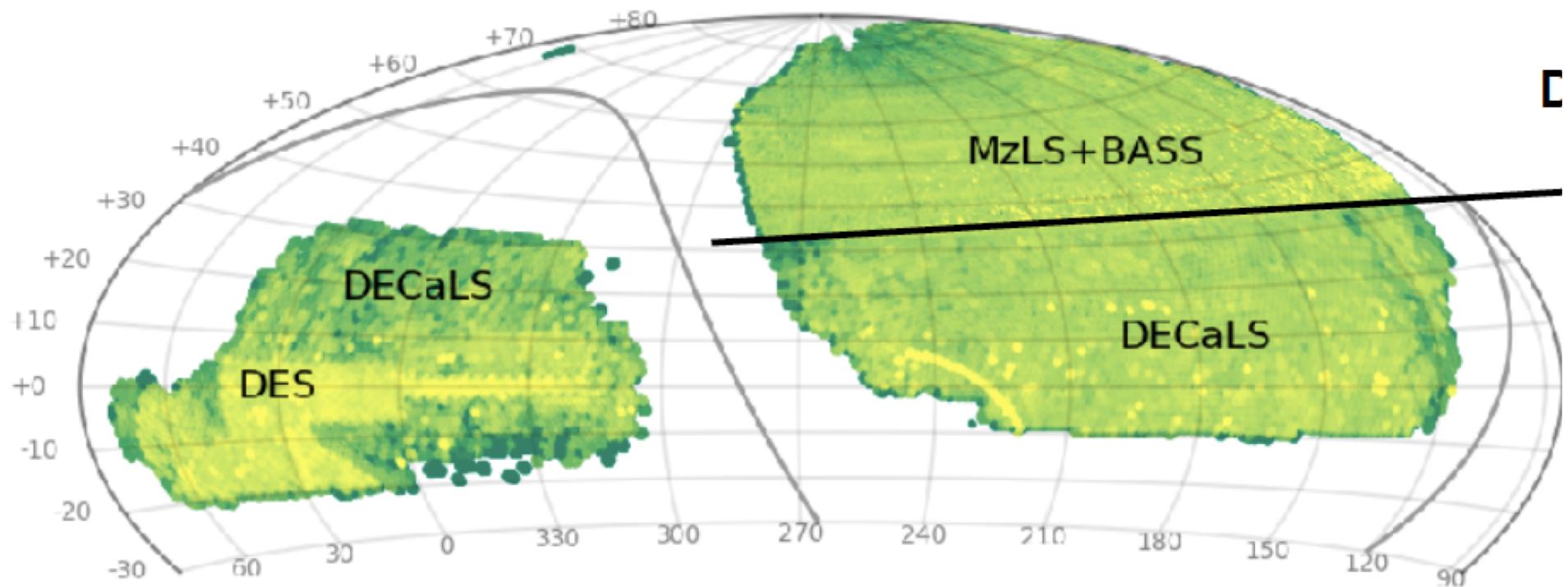
November 05 ,2018





Context

- The Target Selection is based on optical imaging survey developed specially for DESI
- Three optical bands (grz) + two NIR bands (W1,W2)



- QSO selection requires PSF objects
- Stars - QSOs confusion, specially for high-z QSOs ($z > 2.1$)
- Density: 260 deg^{-2} targets



Color cut selection with DECaLS

$$\text{grz flux} = (\text{gflux} + 0.8 \cdot \text{rflux} + 0.5 \cdot \text{zflux}) / 2.4$$

$$\text{W flux} = 0.75 \cdot \text{w1flux} + 0.25 \cdot \text{w2flux}$$

Depth: $\text{W} > 0 \text{ \&\& } \text{grz} > 0 \text{ \&\& } 17 < \text{grz} \text{ \&\& } \text{r} < 22.7$

Morphology: PSF \&\& $\text{dchisq(PSF)} - \text{dchisq(SIM)} > 40$

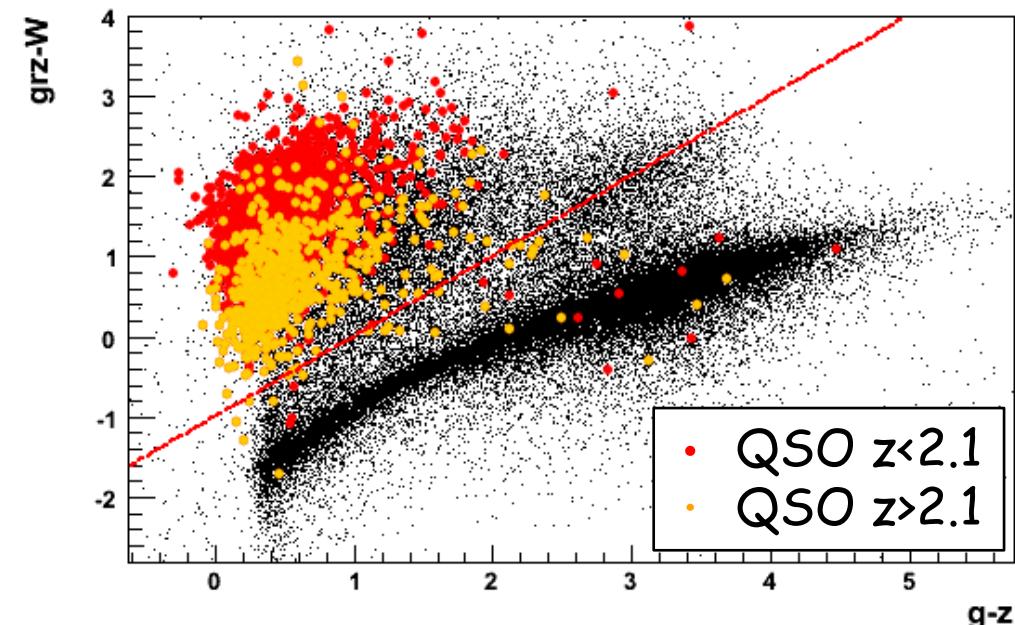
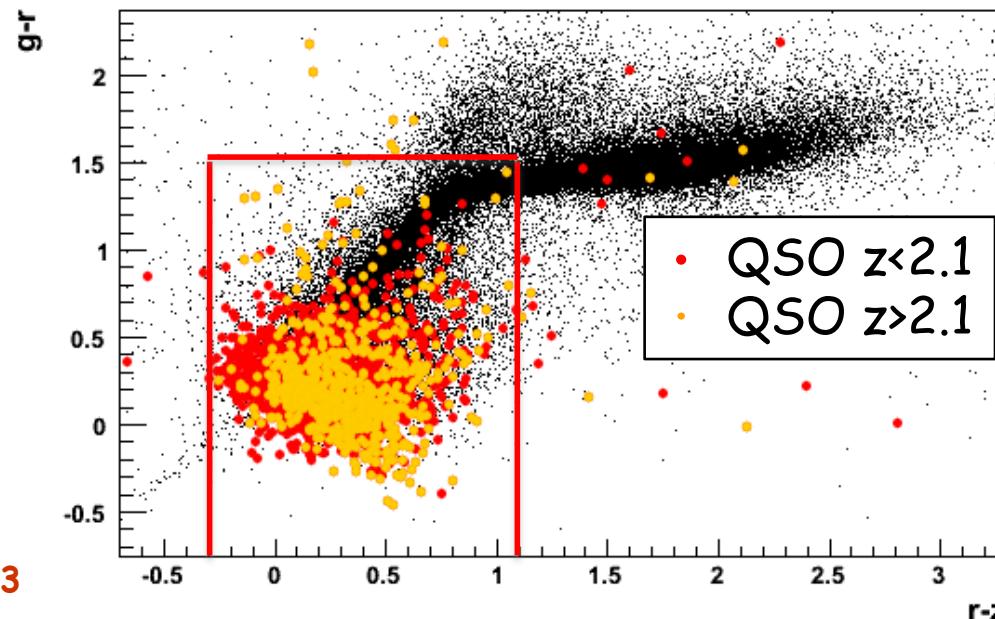
grz box: $(\text{g}-\text{r}) < 1.5 \text{ \&\& } \text{r}-\text{z} < 1.1 \text{ \&\& } \text{r}-\text{z} > -0.3$

grz-W box: $\text{grz-W} > \text{g-z} - 1.0$

W1/W2: $\text{snrW2} > 1 \text{ \&\& } \text{snrW1} > 4 \text{ \&\& } \text{W1-W2} > -0.4$

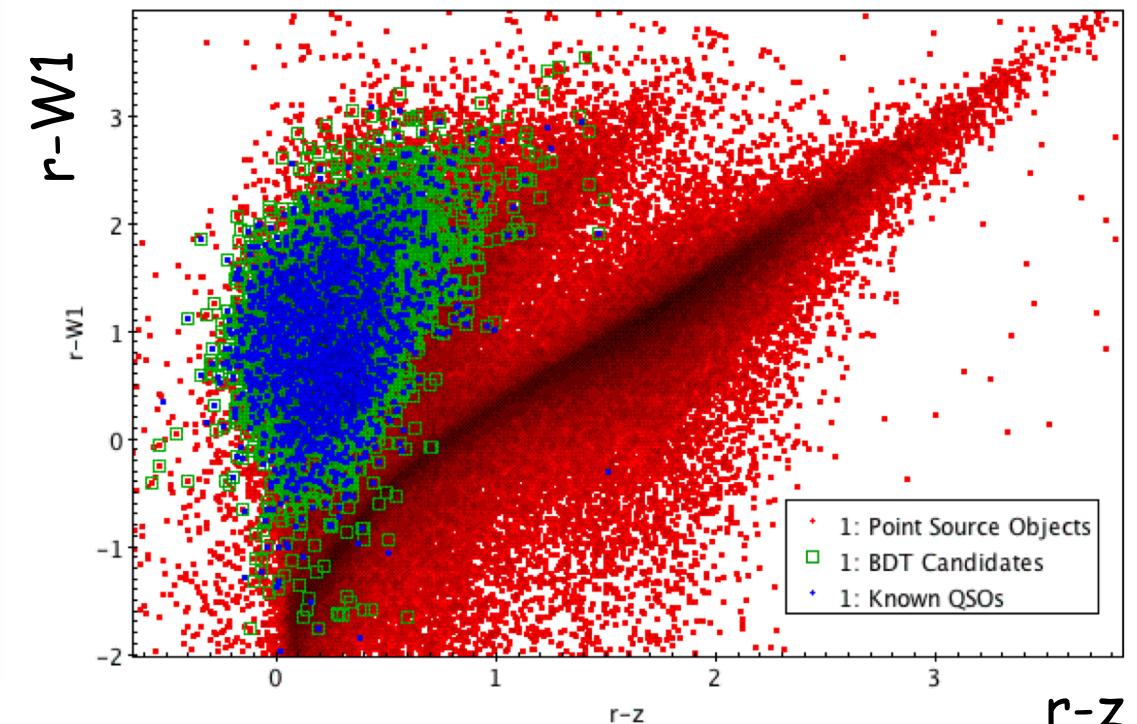
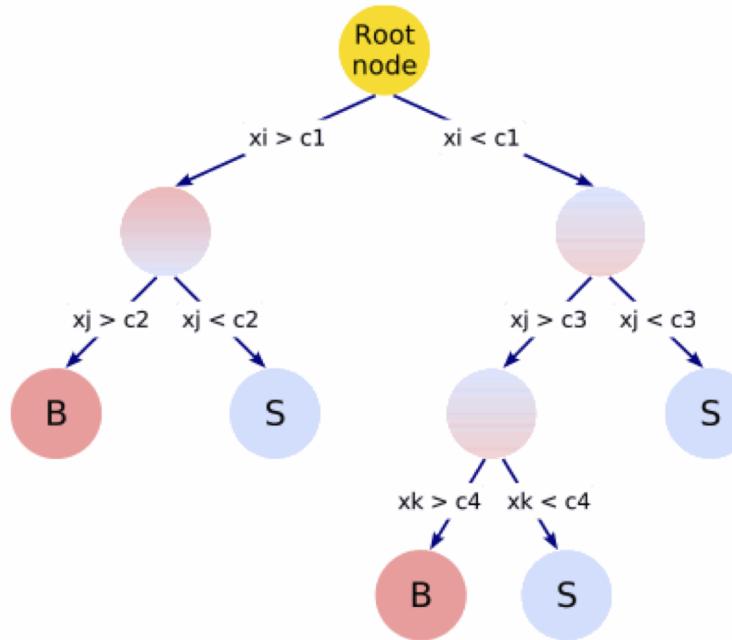
Main sequence rejection: $! (\text{g}-\text{r} > 0.20 \text{ \&\& } \text{abs}((\text{g}-\text{r}) - 1.5 * (\text{r}-\text{z}) - 0.175) < 0.100 \text{ \&\& } (\text{W1-W2} < 0.3 \text{ || } \text{grz-W} > 3))$

What we want to avoid!!!





Machine Learning Approach



Random forest algorithm

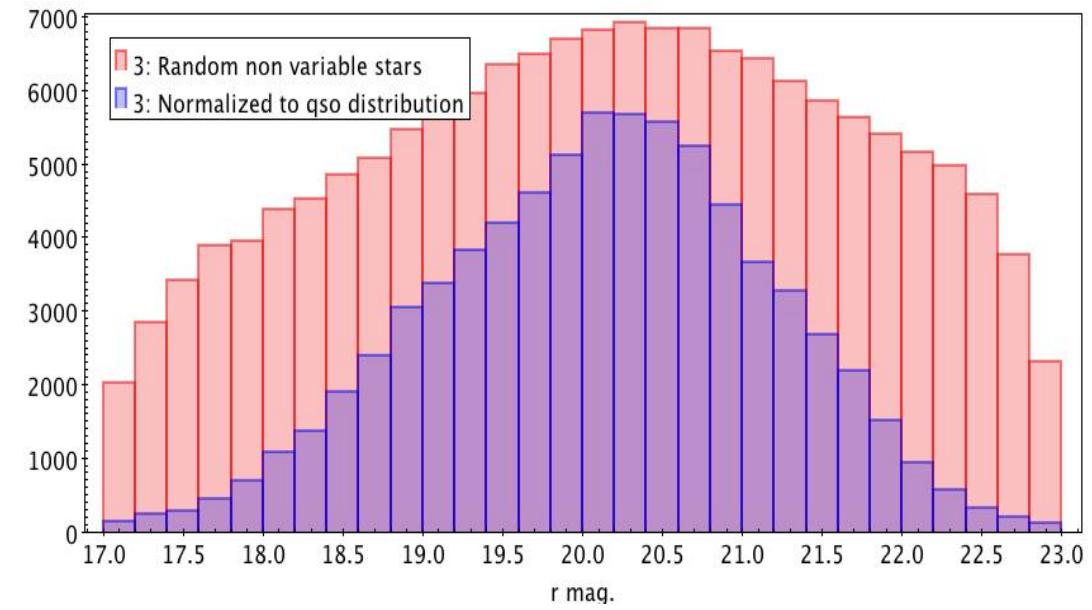
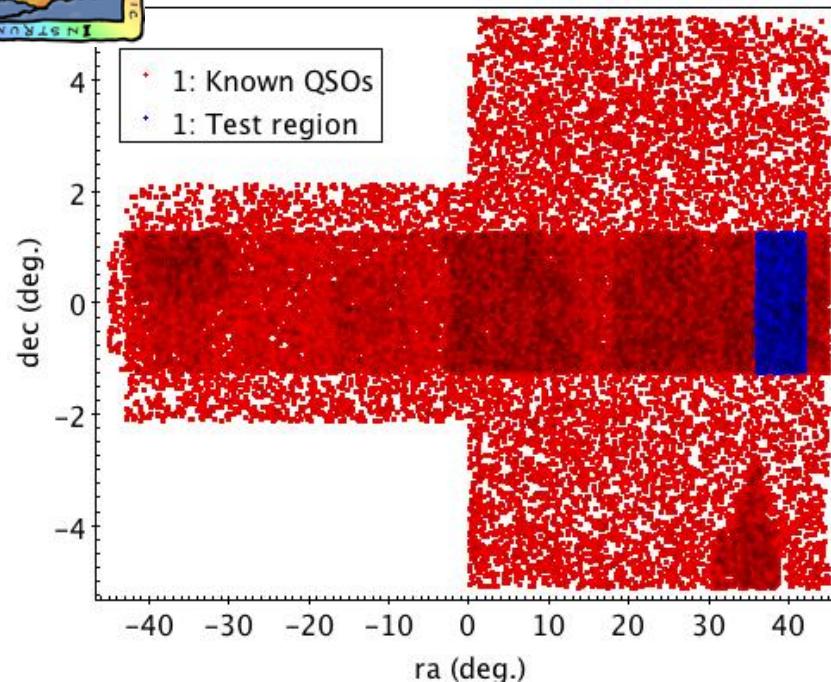
- Classification with **scikit-learn**
- Random: Subset of variables and subsample of the training sample
- Forest: Sum of the decision trees
- Training over all the known QSOs with DECaLS DR3
- Tested in stripe 82

Variables (inputs)

- All possible colors used with DECaLS (g, r, z) and WISE ($W1, W2$)
- a mag. reference: r



Proof of Principle with DR3

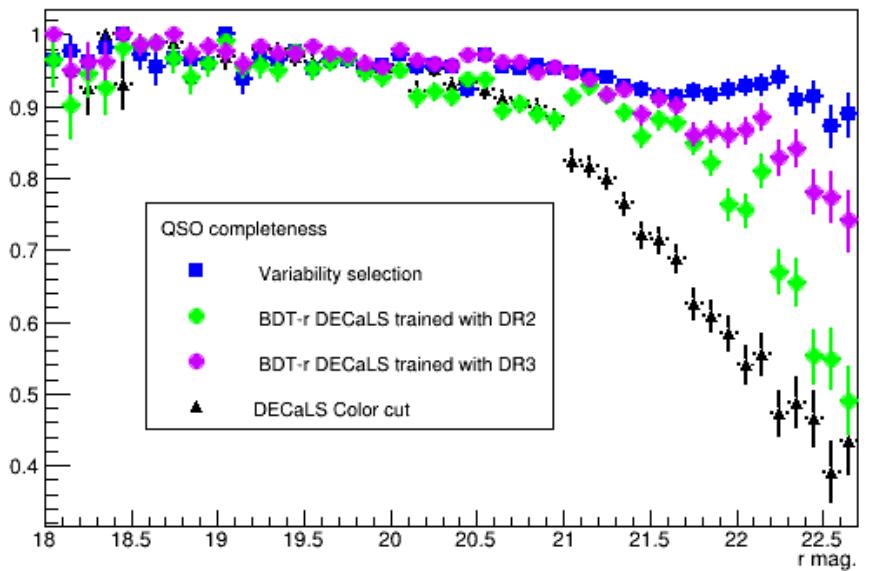
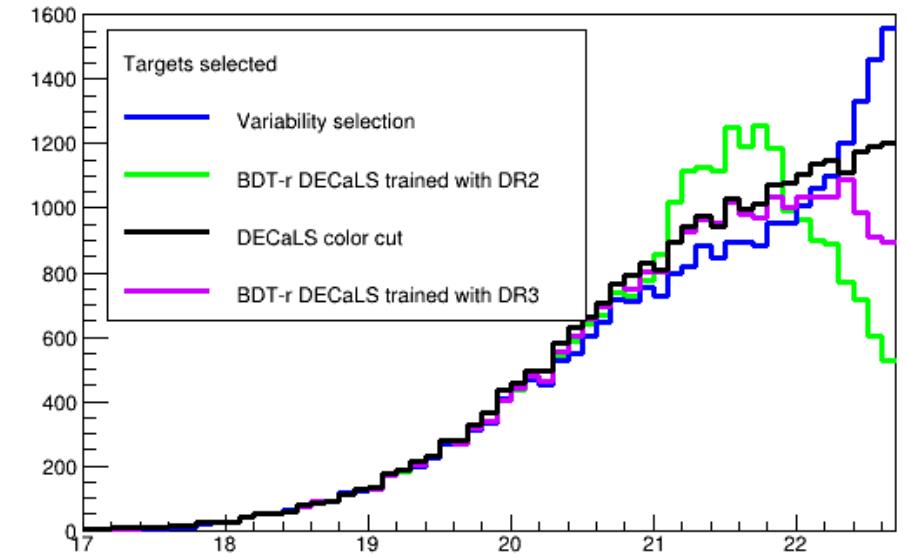


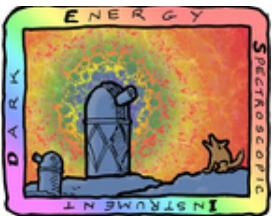
- DR3 DECaLS QSO sample:
 - Deep region of DES (fat stripe 82) : **~35000 objects**
 - Bright objects ($\sigma(r) < 0.02$): **~45000 objects**
- DR3 DECaLS star sample:
 - From fat stripe 82 (known QSOs and variable objects are removed): **~75000 objects**
- To avoid bias, same r mag. distribution for stars and qso (difference between our DR2 and DR3 RF)



Results with Stripe 82

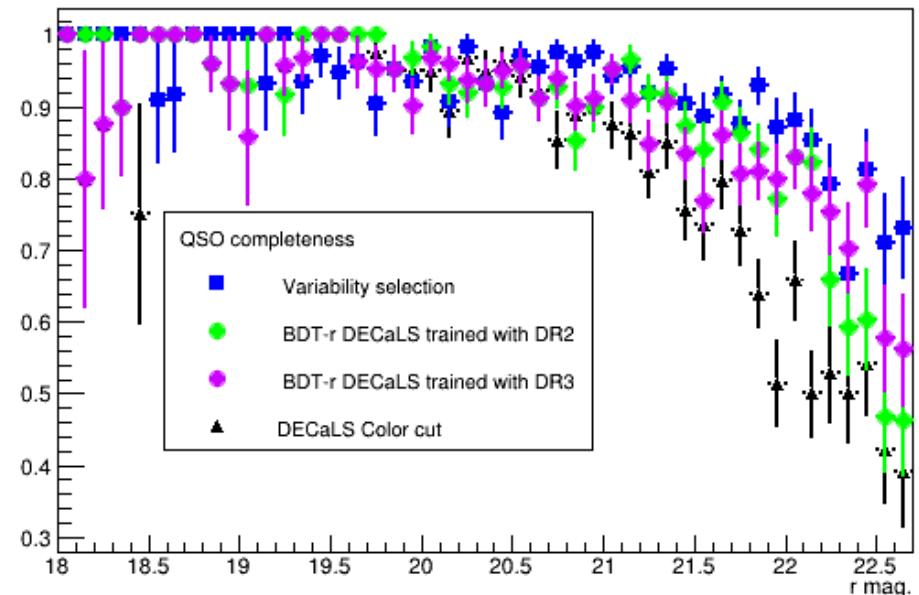
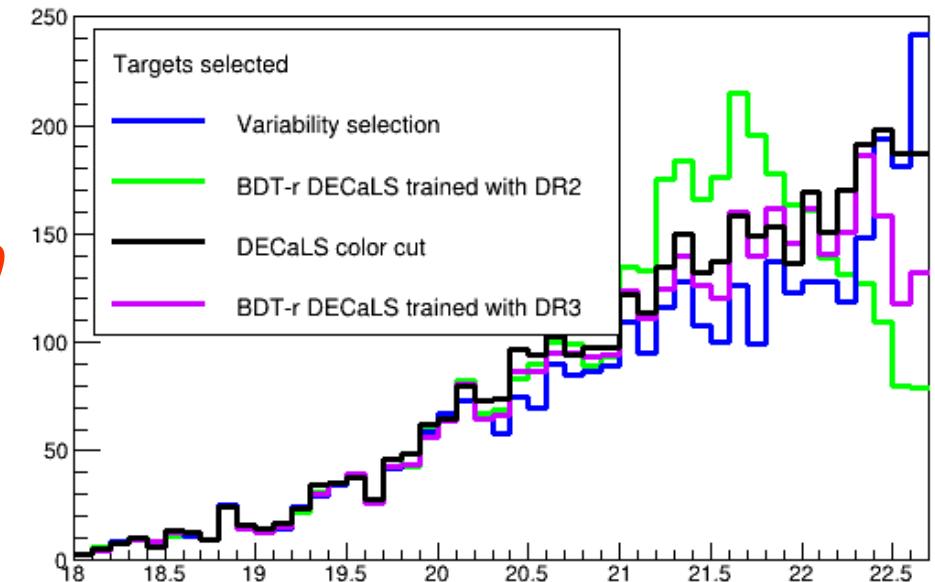
- Test over the stripe 82, $\text{ra} > 0$
- Same budget for the four selections: $\sim 240 \text{ deg}^2$ objects
- Comparison of the Color Cut selection (black triangles) and the RF (purple circles)
- Gain in completeness for faint objects $r > 22$





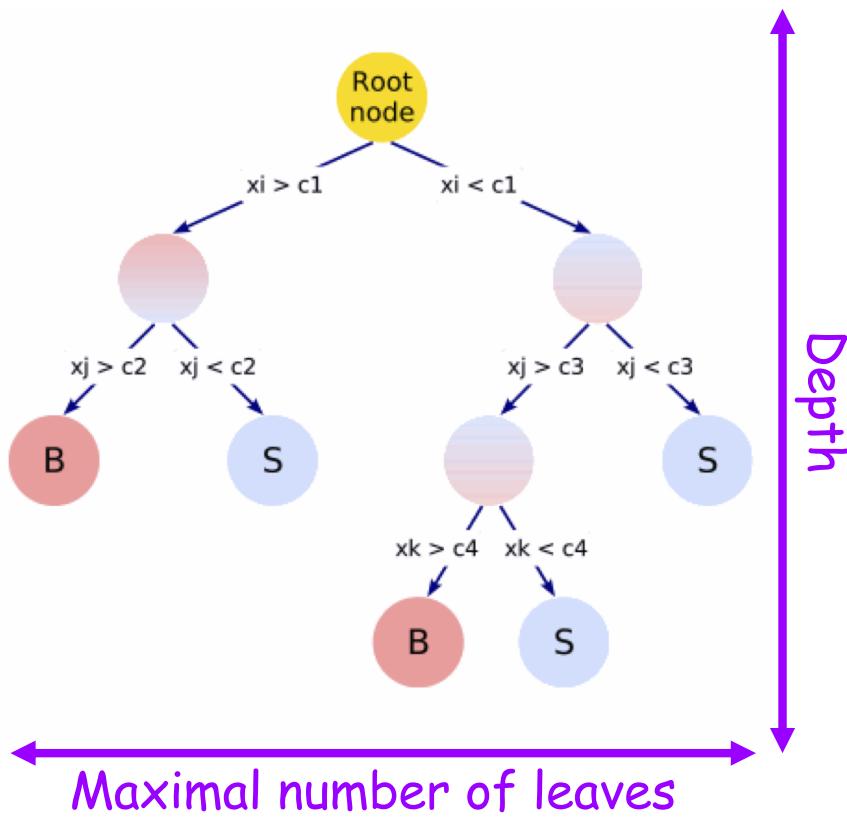
Results with test region in Stripe 82

- Test over the test region of stripe 82, $36 < \text{ra} < 42$
- Same budget : $\sim 265 \text{ deg}^{-2}$
- Comparison of the Color Cut selection (black triangles) and the RF (purple circles) and the Variability selection (blue squares)
- Same trends: Gain in completeness for faint objects $r > 22\dots$





Optimization of the RF Training



Benefits

- Reduce by a factor 3 the size of the RF
- Faster computation of the probability
- Better control of the overtraining

Parameters to avoid overtraining and oversized RF

- Number of trees in the forest
- The maximum depth of the tree
- The maximum number of leaves

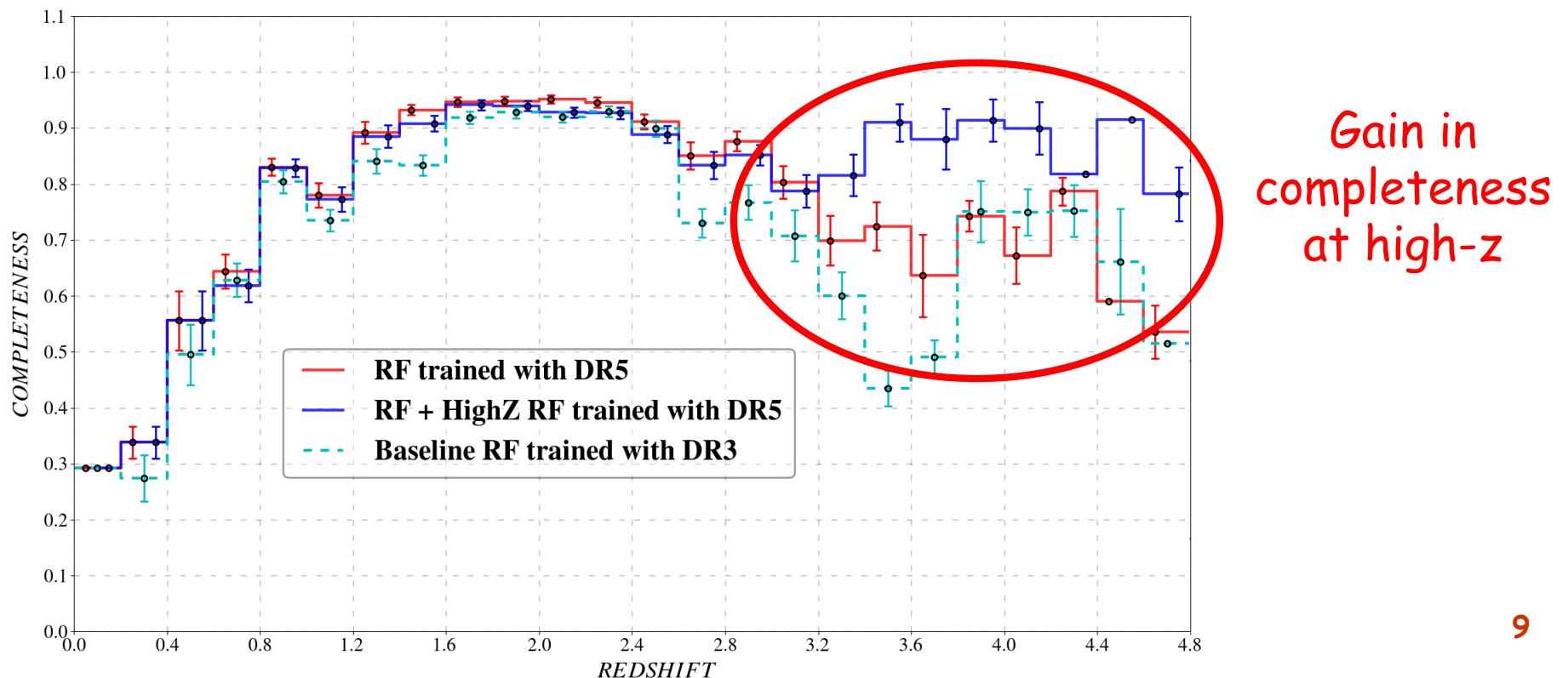
Scikit-learn implementation

- **RandomForestClassifier**
- `n_estimators`
- `max_depth`
- `max_leaf_nodes`
- As a result the last node (leaf) is not pure in Signal/Background



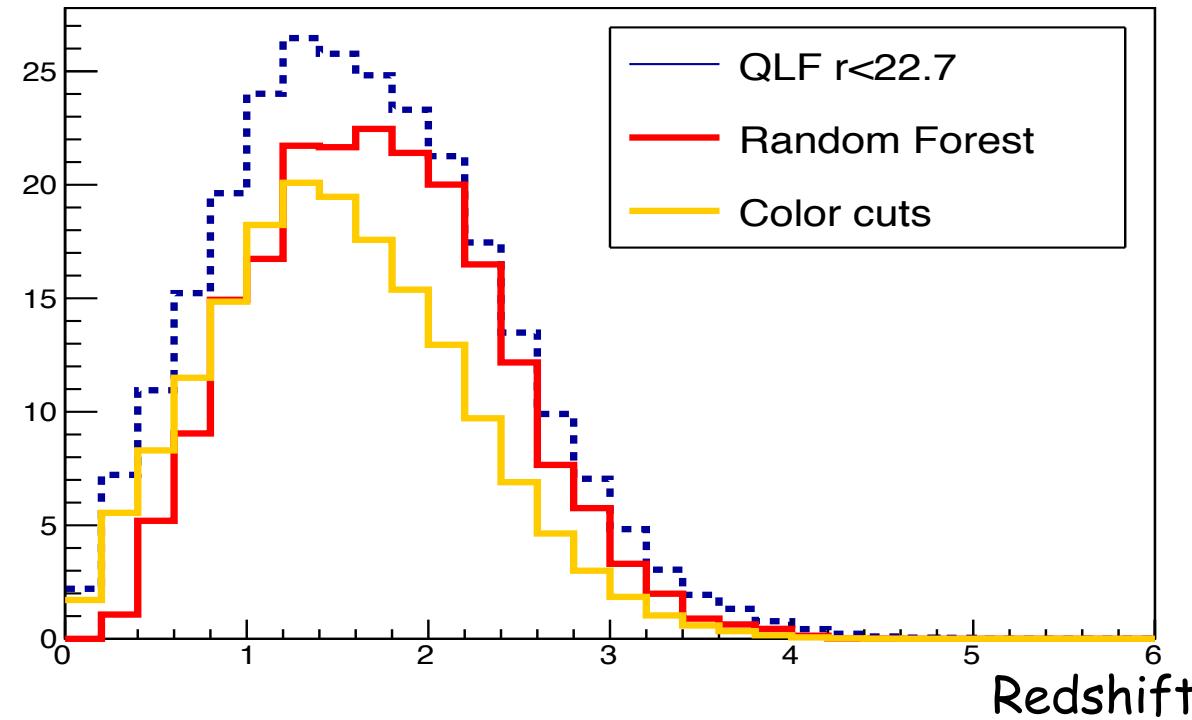
Latest developments

- Two different RF:
 - RF trained for all redshifts, 245 deg^{-2}
 - RF trained for $z > \sim 3.5$, 15 deg^{-2}
- Different training for DR3, DR5, DR7
 - Very stable results (probability cut slightly tuned)





Conclusions



- Very simple to implement in scikit-learn
- Own persistency of the RF in *desitarget* code
- Very stable as a function of the various imaging releases (DR5, DR7...)
- Easy to tune to achieve the target budget
- Much efficient than the color cut approach specially for high-z QSOs