

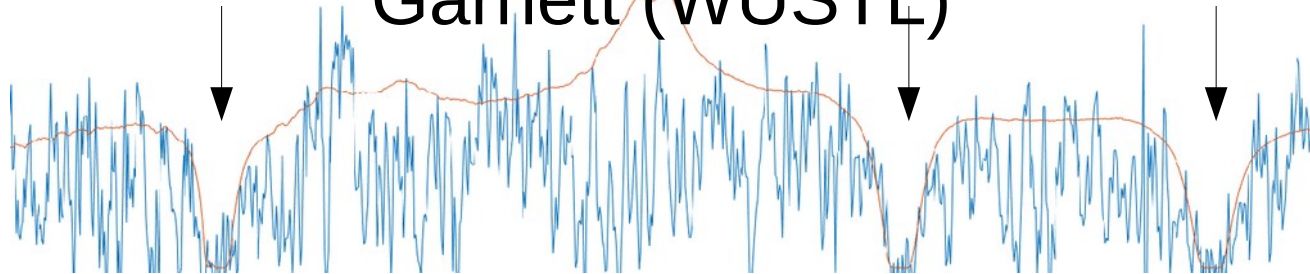


# Finding Strong Ly $\alpha$ Absorbers in the Shadows of Quasars with Bayesian Machine Learning

 <https://arxiv.org/abs/2003.11036>

 [https://github.com/rmgarnett/gp\\_dla\\_detection/](https://github.com/rmgarnett/gp_dla_detection/)

Ming-Feng Ho (Me) (UCR), Simeon Bird (UCR), Roman  
Garnett (WUSTL)



**DLAs (Damped Lyman alpha absorbers):**

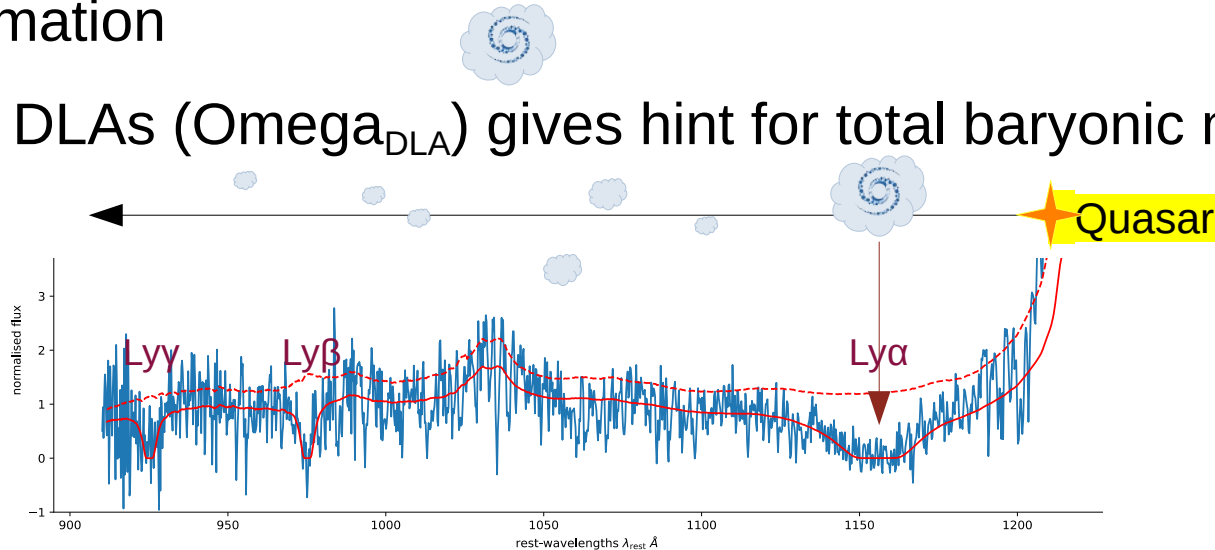
Strong neutral hydrogen absorbers (usually  $2 < z < 5$ ).

Dominate neutral hydrogen budget after reionisation.



# What are Damped Lyman alpha absorbers (DLAs)?

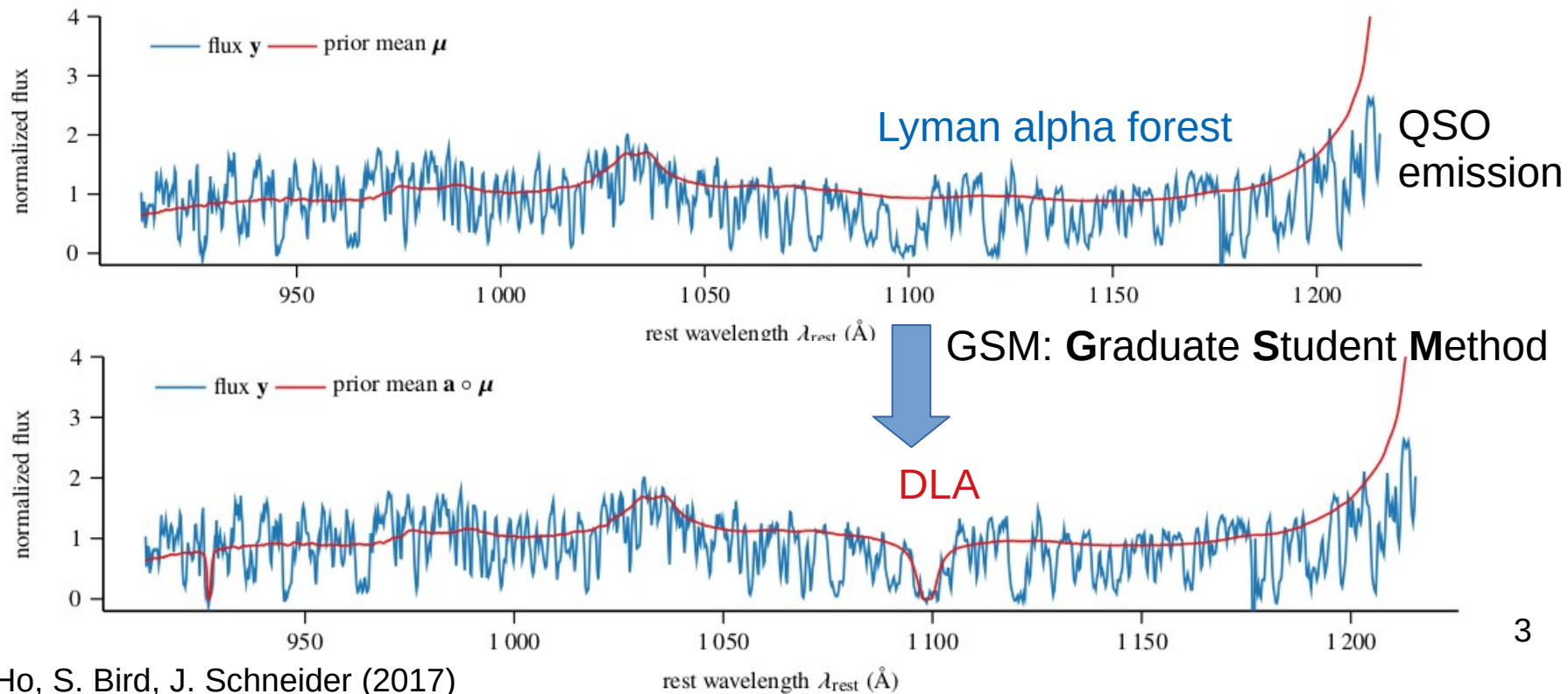
- Neutral hydrogen gas with a **high column density** ( $>10^{20.3} \text{ cm}^{-2}$ )
- **Baryonic acoustic oscillation** (BAO): DLAs, uncertainty in Lyman alpha forest power spectrum
- Ultimately accretes onto galactic halos and **fuels star formation**: hint for galaxy formation
- Total mass of DLAs ( $\Omega_{\text{DLA}}$ ) gives hint for total baryonic matter ( $\Omega_{\text{b}}$ )



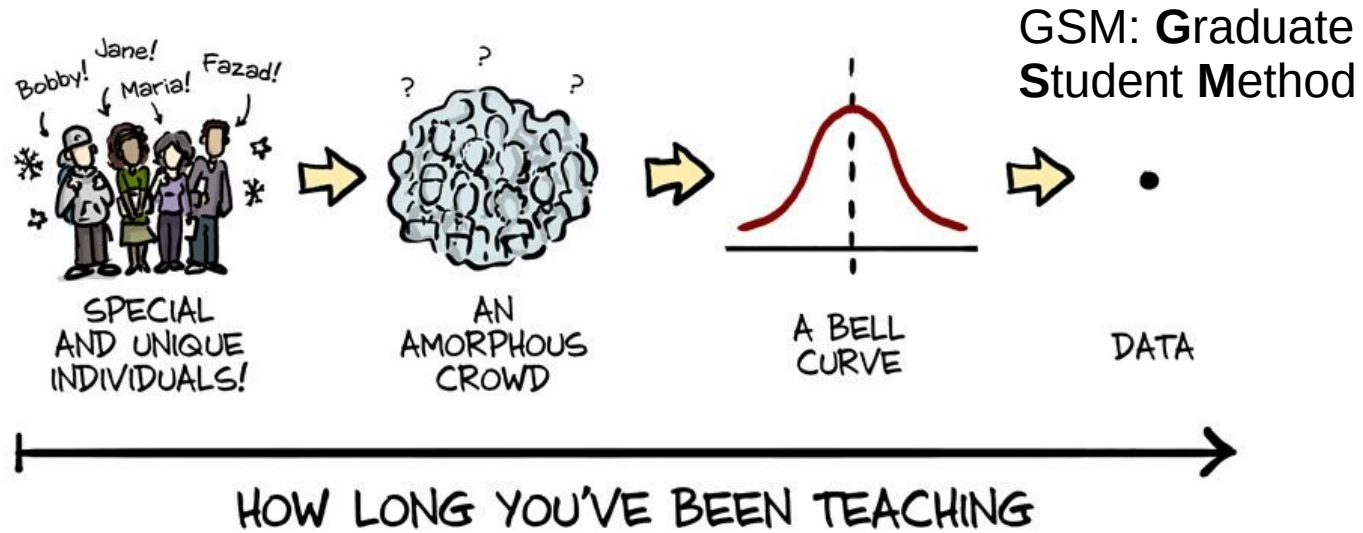
# Finding DLAs in Spectra

Currently done by **visual inspection** of spectra

Look for wide dips in the spectrum below (through GSM):

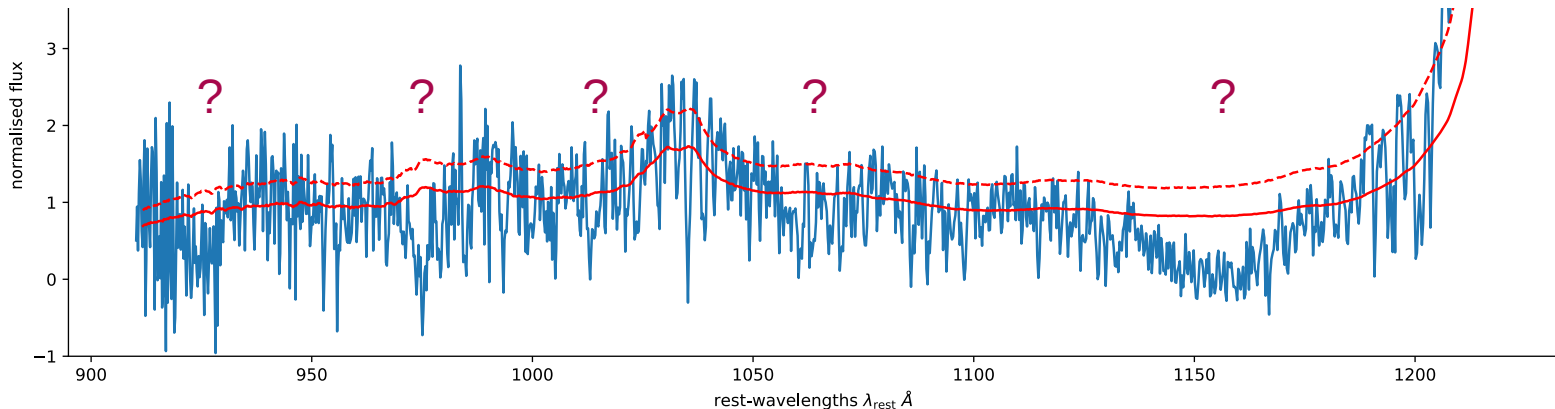


# HOW YOU SEE YOUR STUDENTS:



# Why Machine Learning?

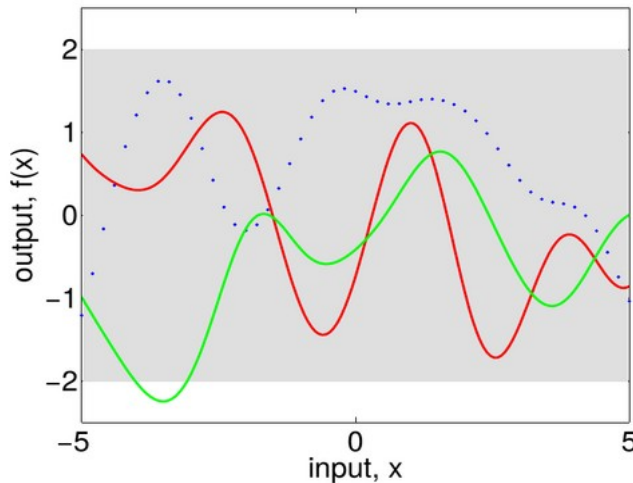
- State-of-art: **visual inspection**, graduate student **method** (GSM)
- **No physical model** for quasar emission yet
- Finding DLAs out of weak absorbers in the forest is hard



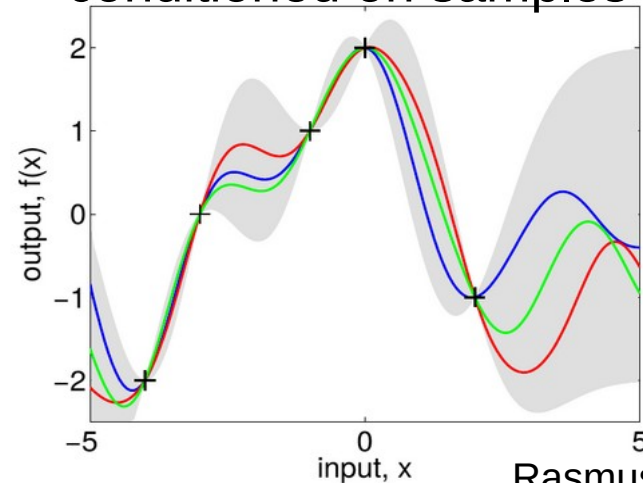
# What is a Gaussian Process?

- Quasar model is a Gaussian Process
- Bayesian function interpolation, which computes probability distribution of  $f(x)$  conditional on input set.

Left: Prior on function



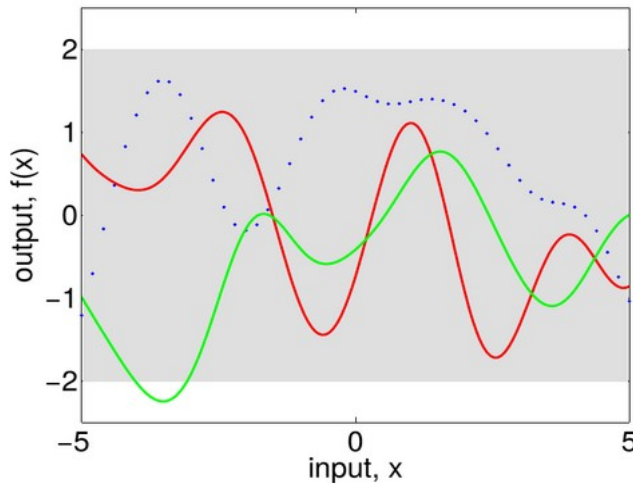
Right: Posterior conditioned on samples



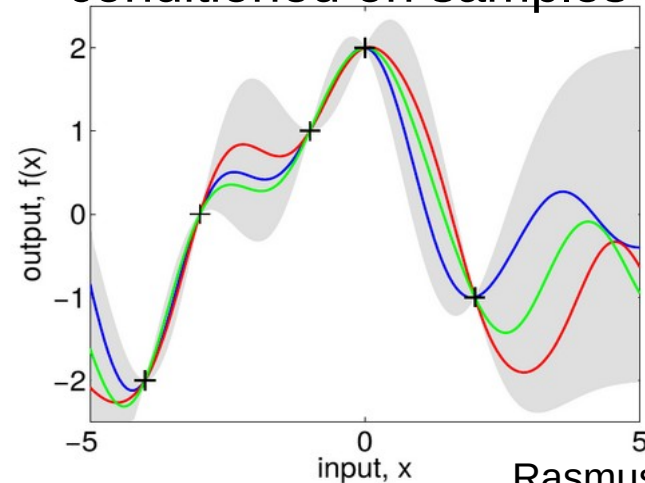
# What is a Gaussian Process?

- Magic in **Kernel** function: describes how **correlation between function values** depends on parameter distance.
- We trained our kernel from quasar spectra; it describes the correlation between different emission lines.

Left: Prior on function



Right: Posterior conditioned on samples



# Bayesian Model Selection

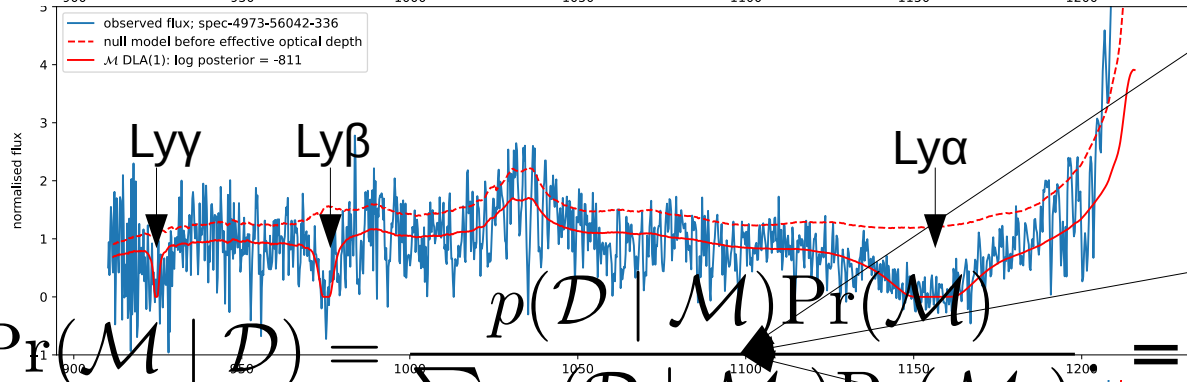
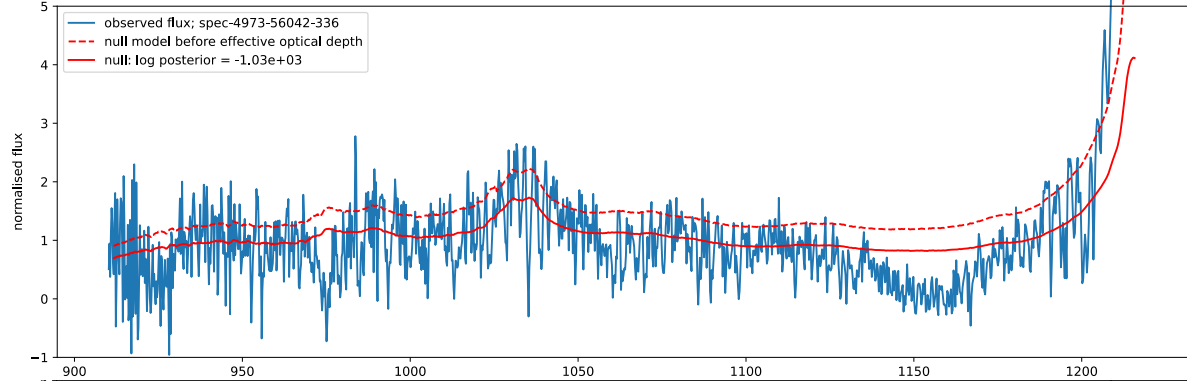
- Trick: train a GP on spectra **without DLAs**
- Build another GP for spectra **with DLAs**
- Evaluate **model posterior**:

$$\Pr(\mathcal{M} \mid \mathcal{D}) = \frac{p(\mathcal{D} \mid \mathcal{M})\Pr(\mathcal{M})}{\sum_i p(\mathcal{D} \mid \mathcal{M}_i)\Pr(\mathcal{M}_i)}.$$



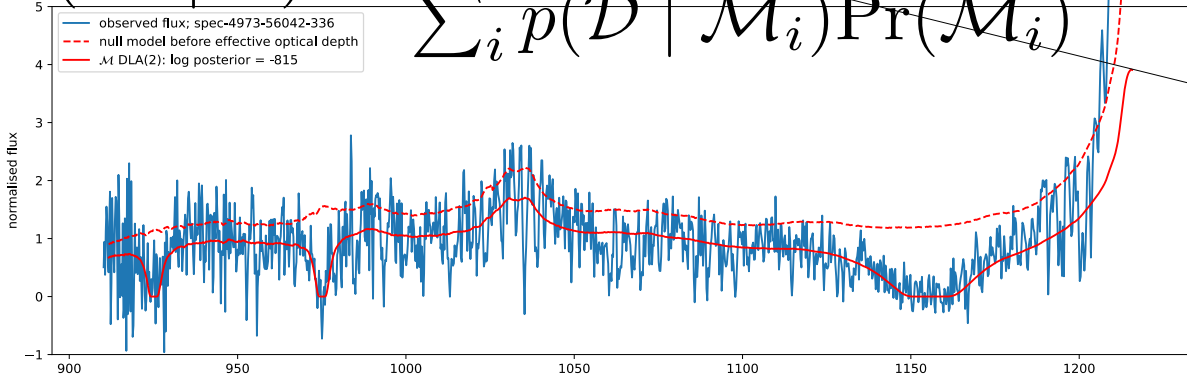
# Bayesian Model Selection

$\log P(0 \text{ DLA} | D) = -1030$



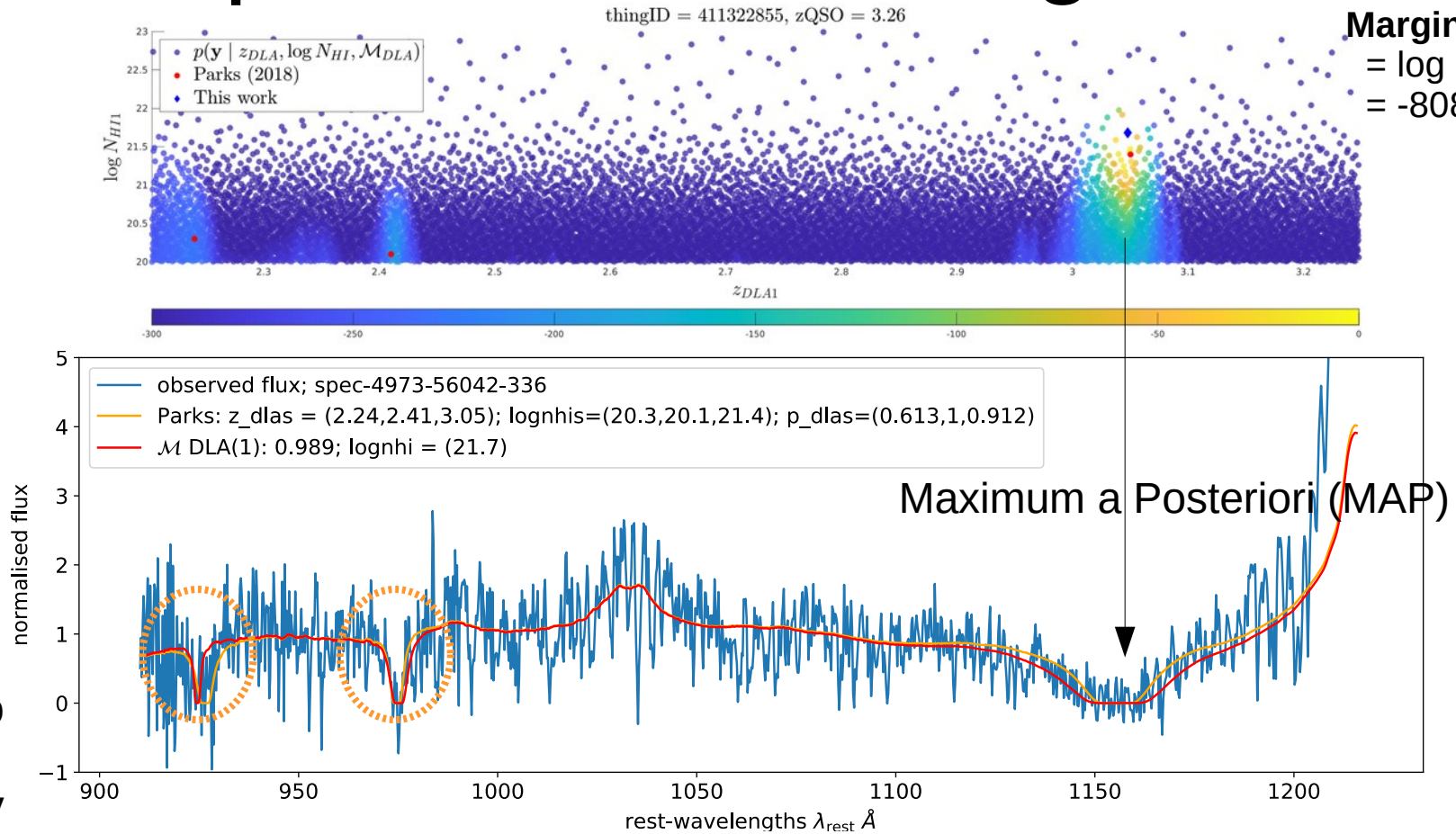
**Select!**  
 $\log P(1 \text{ DLA} | D) = -811$

$$\Pr(\mathcal{M} | D) = \frac{p(D | \mathcal{M}) \Pr(\mathcal{M})}{\sum_i p(D | \mathcal{M}_i) \Pr(\mathcal{M}_i)} = 0.989$$



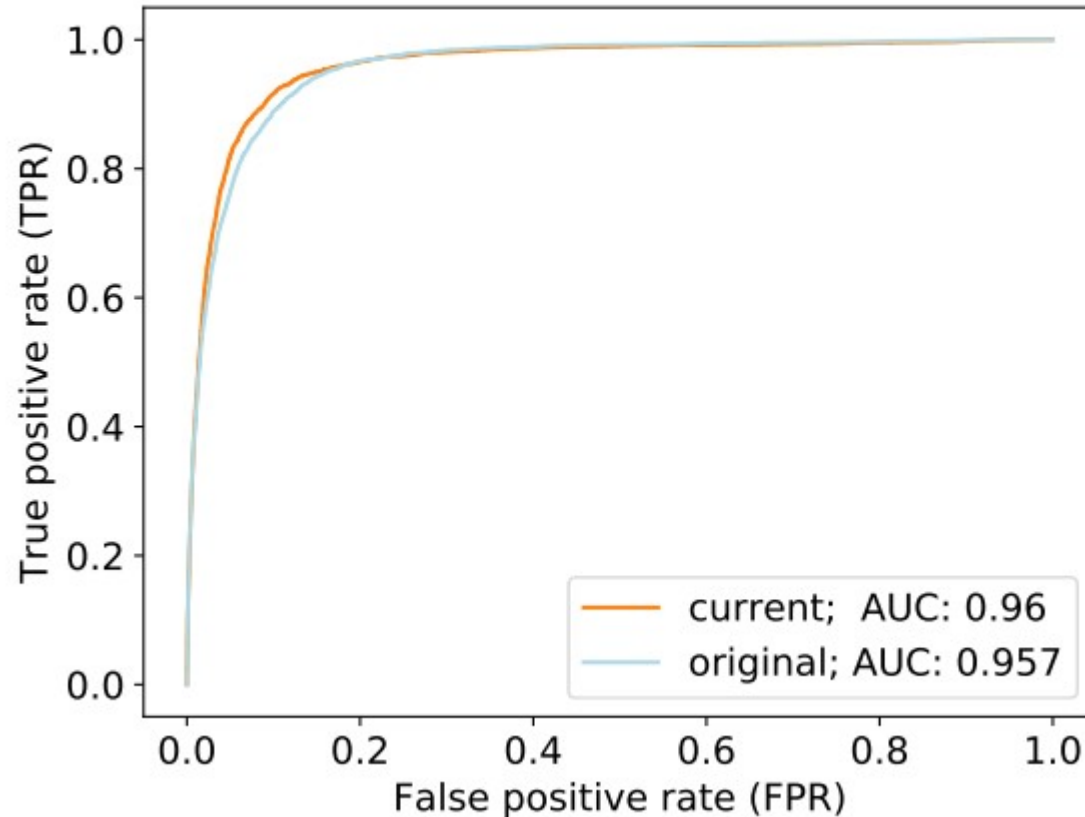
$\log P(2 \text{ DLAs} | D) = -815$

# Sum up all possible DLAs $\rightarrow$ the posterior of having DLAs



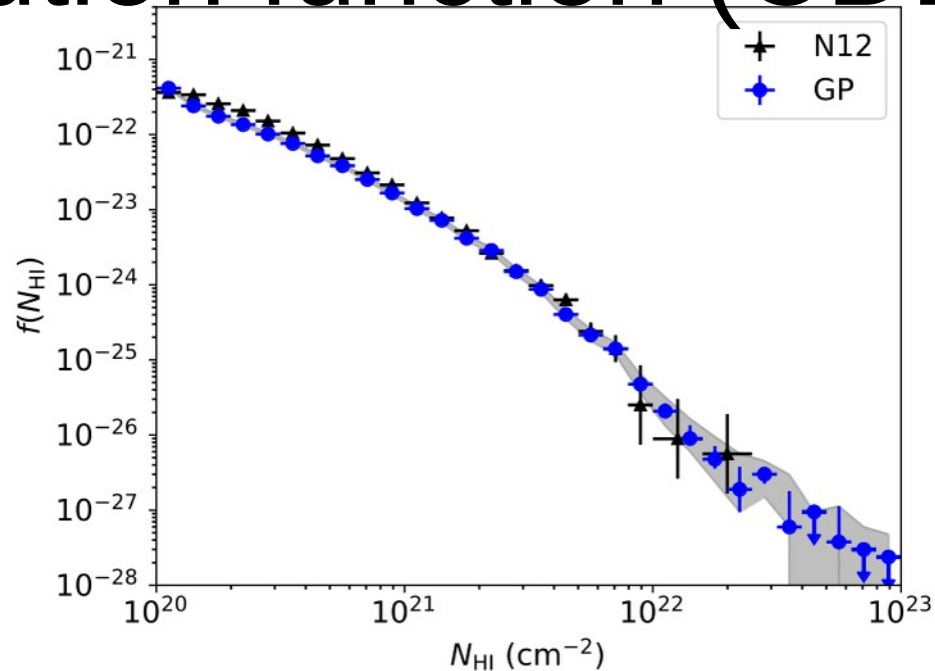
Important to  
also model  
 $\text{Ly}\beta$  and  $\text{Ly}\gamma$

# Results: A decent ROC



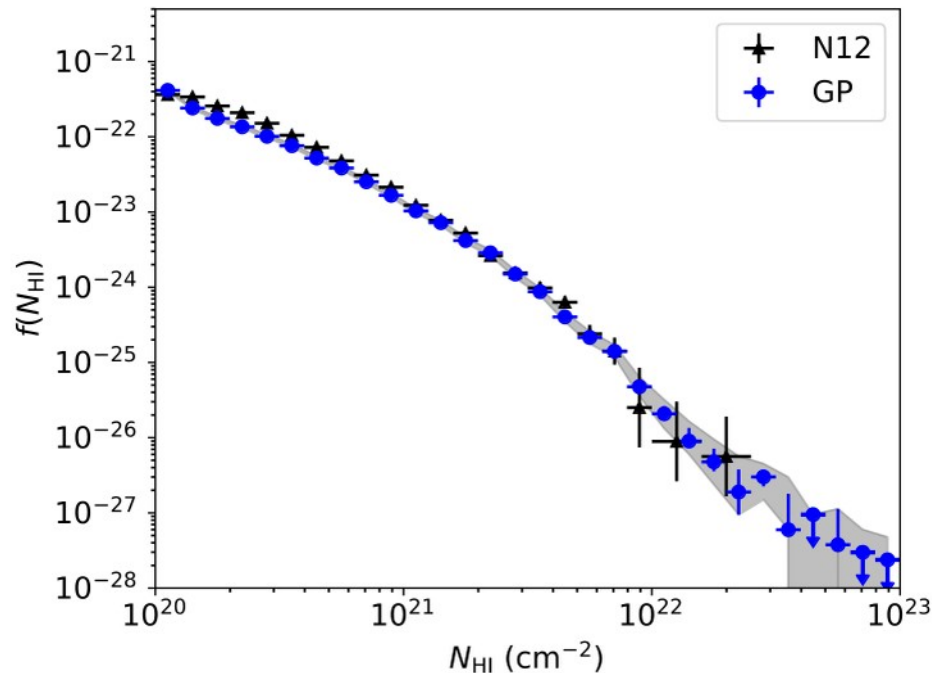
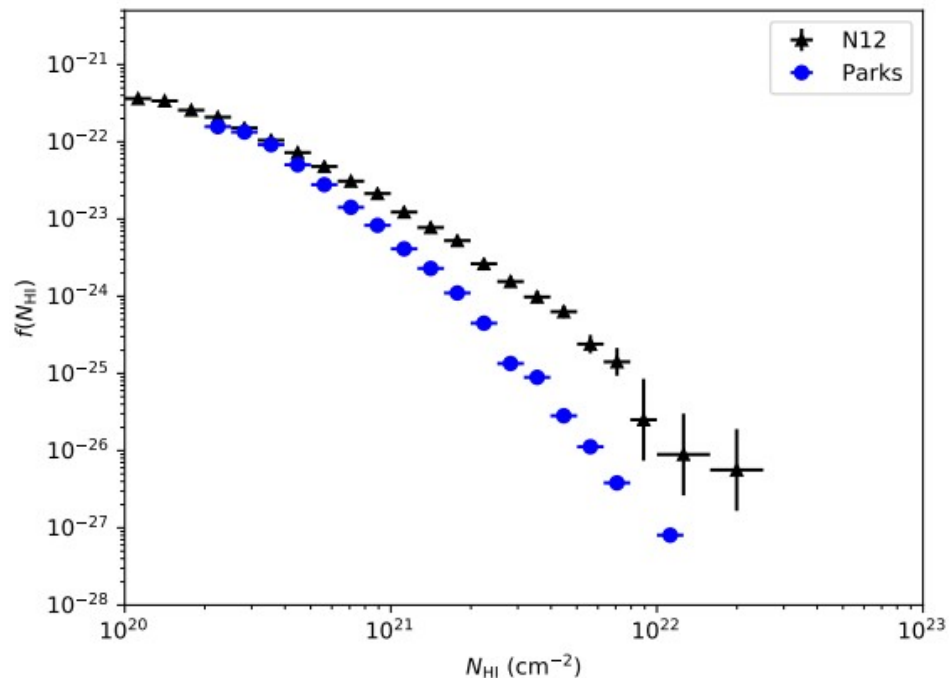
A better handle for those borderline cases with  $\log\text{NHI} \sim 20$

# Results: Column density distribution function (CDDF)



- Use all data (DR12), even with  $\text{SNR} < 1$
- Properly Propagate the uncertainty to  $N_{\text{HI}}$
- Reproduce previous catalogue (Noterdaeme 2012)

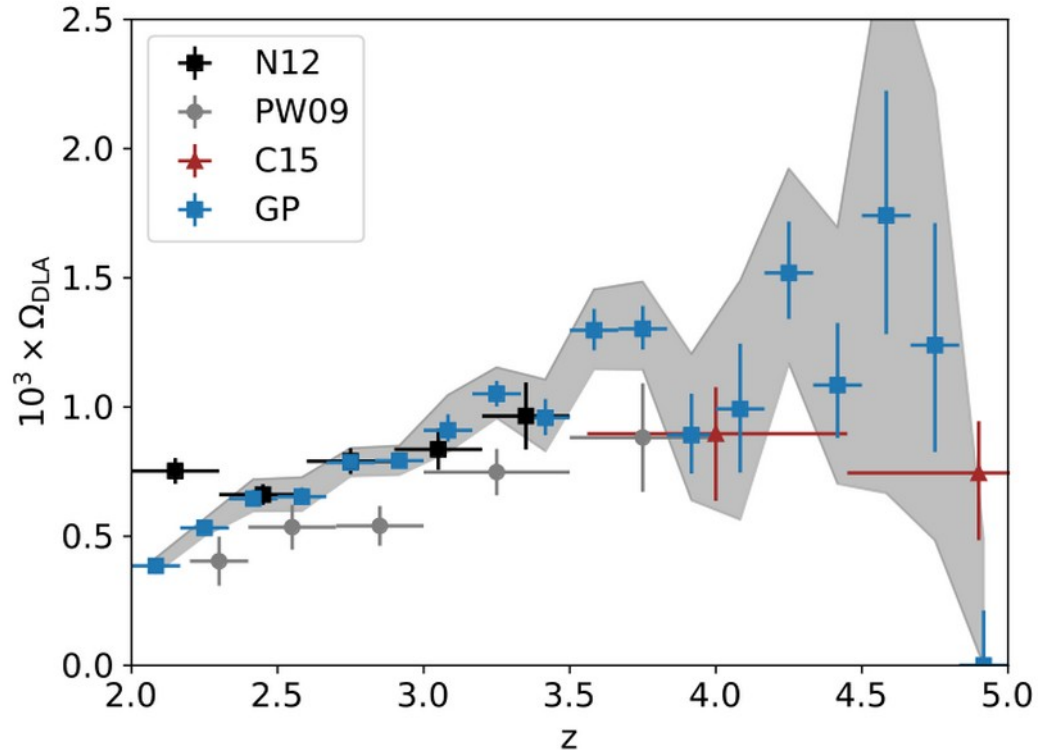
# Parks' CDDF



Note: this CDDF is reproduced by us using Parks' catalogue.

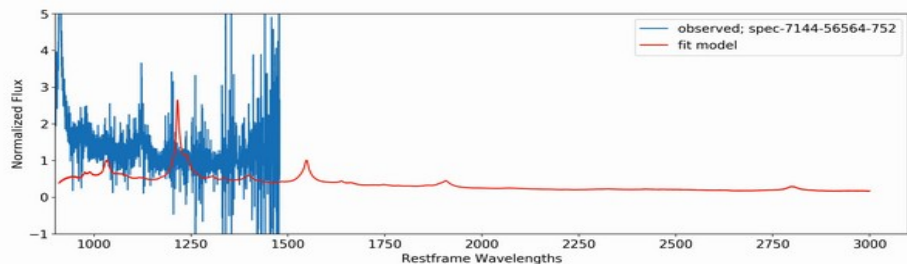
- More precise measurement than Parks on the **high-end**

# Results: Total Mass of DLA

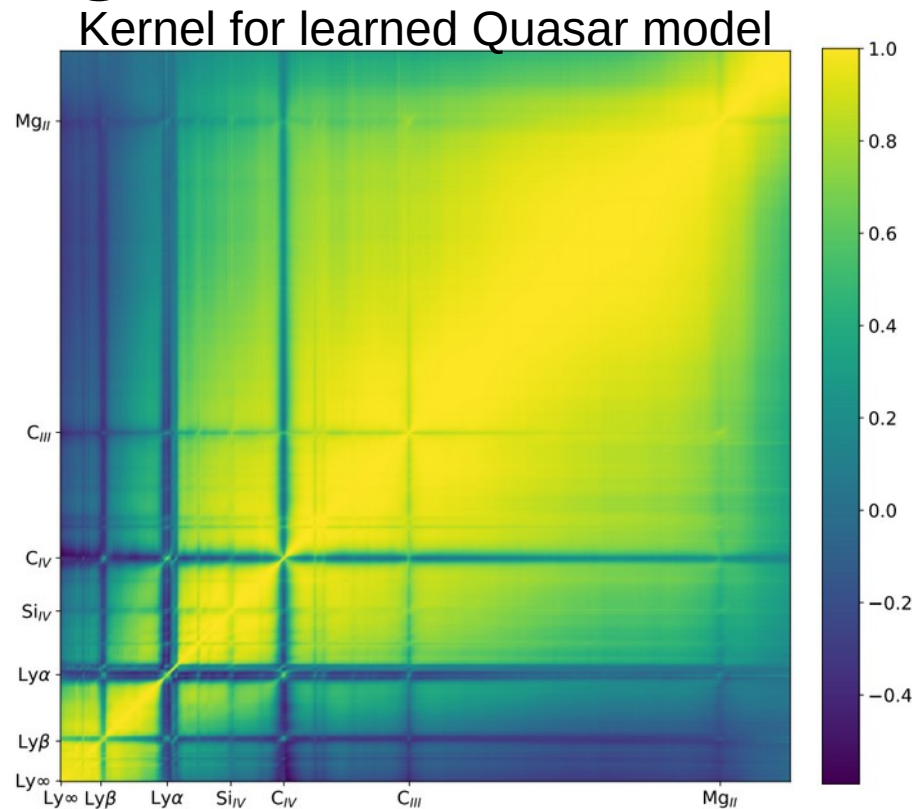


- Count all **weak detections**
- Extend measurements to  **$z > 4.0$**

# Bonus: Quasar redshift estimation using a GP




with UCR CS people:  
Leah Fauber, Christian Shelton, Ishita Korde




# Conclusion

- **Automated** the detection of **DLAs**
- We get a **posterior density** per spectrum
- Decent accuracy; better NHI estimations
- Can also **estimate zQSO**

 <https://arxiv.org/abs/2003.11036>

 Our model is publicly available :  
[https://github.com/rmgarnett/gp\\_dla\\_detection/](https://github.com/rmgarnett/gp_dla_detection/)

 Me:  
<https://github.com/jibanCat>

 Redshift estimation code with GP:  
[https://github.com/sbird/gp\\_qso\\_redshift](https://github.com/sbird/gp_qso_redshift)

Thank Reza Mondai & Madi Qezlou(UCR Astro):  
for valuable comments.

Thank Yongda Zhu and Marie Wingyee Lau (UCR, Astro):  
for useful discussions on quasar continua.

Thank Julia Plank (UAuckland, Psych):  
for listening to my complaints oversea psychological supports during pandemic



How do we count expected number of DLAs at a given redshift bin?

# How we count DLAs on column density bins?

- Having posterior density: larger sample size and smaller error bars
- Compute  $\Pr(N)$  with Poisson-Binomial process

$$\Pr(N) = \sum_{\text{DLA} \in F_N} \prod_{i \in \text{DLA}} p_{\text{DLA}}^i(\{\mathcal{M}_{\text{DLA}}\} | Q) \prod_{j \in \text{DLA}^c} (1 - p_{\text{DLA}}^j(\{\mathcal{M}_{\text{DLA}}\} | Q))$$

where  $F_N$  corresponds to all subsets of  $N$  integers can be selected from the sequence  $\{1, 2, \dots, n\}$ . Bin  $Q$  is an interval in the parameter space of column density or DLA redshift  $Q \in \{N_{\text{HI}}, z_{\text{DLA}}\}$