

Extreme value statistics in astronomy

David Valls-Gabaud

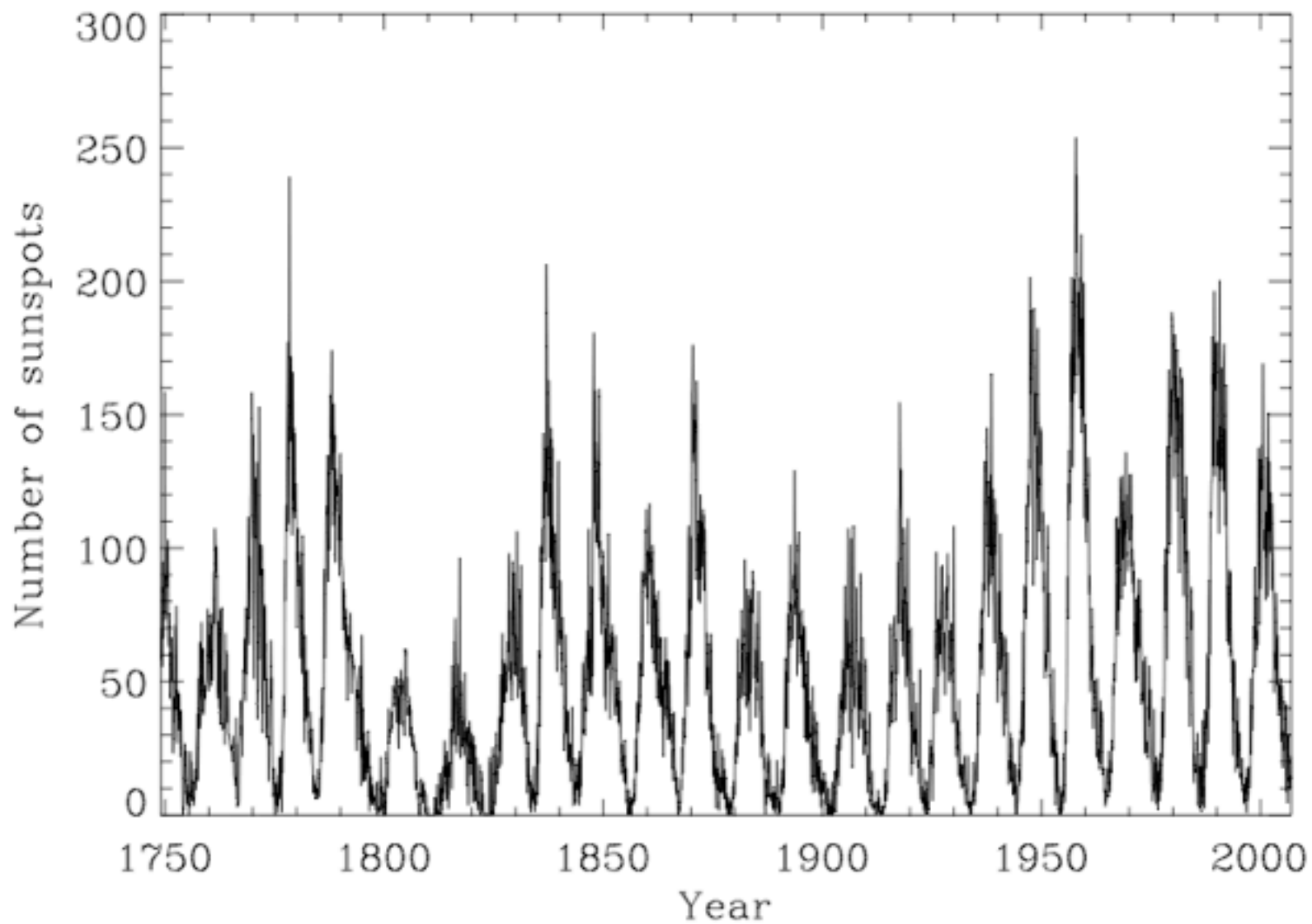
CNRS, Observatoire de Paris



Atelier Astrostatistiques
Grenoble 2011 Dec 09



Extreme values of a sample



Basics of extreme value theory I

Let $X_1 \dots X_n$ be independent random variables having a common distribution function F .

Let M_n be the maximum of n observations :

$$M_n = \max \{X_1, \dots, X_n\}$$

The theoretical distribution of M_n

$$\begin{aligned} \Pr \{M_n \leq z\} &= \Pr \{X_1 \leq z, \dots, X_n \leq z\} \\ &= \Pr \{X_1 \leq z\} \times \dots \times P \{X_n \leq z\} \\ &= \{F(z)\}^n . \end{aligned}$$

not very useful since usually F is unknown...

Basics of extreme value theory II

Make a simple linear renormalisation of M_n

$$M_n^* = \frac{M_n - b_n}{a_n}$$

for constants $\{a_n > 0\}$ and $\{b_n\}$ chosen appropriately to stabilise the location and scale of M_n^* as n increases.

Basics of extreme value theory III

Theorem (Fréchet 1927): If there exist sequences of constants $\{a_n > 0\}$ and $\{b_n\}$ such that

$$\Pr \left\{ \frac{M_n - b_n}{a_n} \leq z \right\} \rightarrow G(z) \quad \text{as } n \rightarrow \infty,$$

where G is a non-degenerate distribution function, then G belongs to one of the following families:

$$\text{I : } G(z) = \exp \left\{ - \exp \left[- \left(\frac{z - b}{a} \right) \right] \right\}, \quad -\infty < z < \infty;$$

$$\text{II : } G(z) = \begin{cases} 0, & z \leq b, \\ \exp \left\{ - \left(\frac{z - b}{a} \right)^{-\alpha} \right\}, & z > b; \end{cases}$$

$$\text{III : } G(z) = \begin{cases} \exp \left\{ - \left[- \left(\frac{z - b}{a} \right)^\alpha \right] \right\}, & z < b, \\ 1, & z \geq b, \end{cases}$$

independently of the underlying F ... (REM: central limit theorem for sample means)

- **The Gumbel distribution**

$$G(x) = \exp(-\exp(-x))$$



Emil Julius Gumbel
1891-1966

- **The Fréchet distribution**

$$G(x) = \begin{cases} 0 & x \leq 0 \\ \exp(-x^{-\alpha}) & x > 0, \alpha > 0 \end{cases}$$



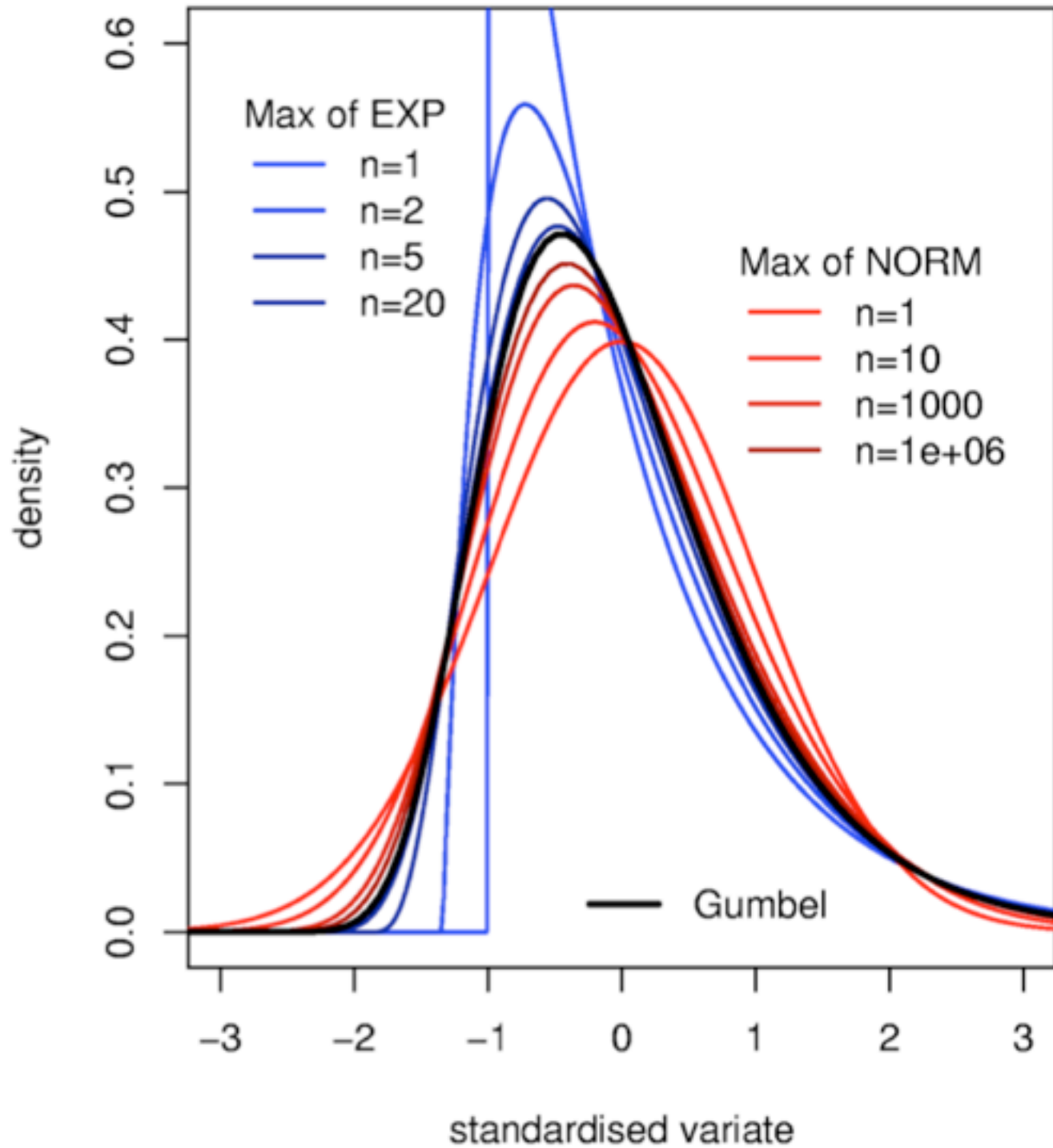
Maurice René
Fréchet
1878-1973

- **The Weibull distribution**

$$G(x) = \begin{cases} \exp(-(-x)^{\alpha}) & x < 0, \alpha > 0 \\ 1 & x \geq 0 \end{cases}$$



Ernst Hjalmar
Waloddi Weibull
1887-1979



Corollary:

If there exist sequences of constants $\{a_n > 0\}$ and $\{b_n\}$ such that

$$\Pr \left\{ \frac{M_n - b_n}{a_n} \leq z \right\} \rightarrow G(z) \quad \text{as } n \rightarrow \infty,$$

for a non-degenerate distribution function G , then G is a member of the GEV family

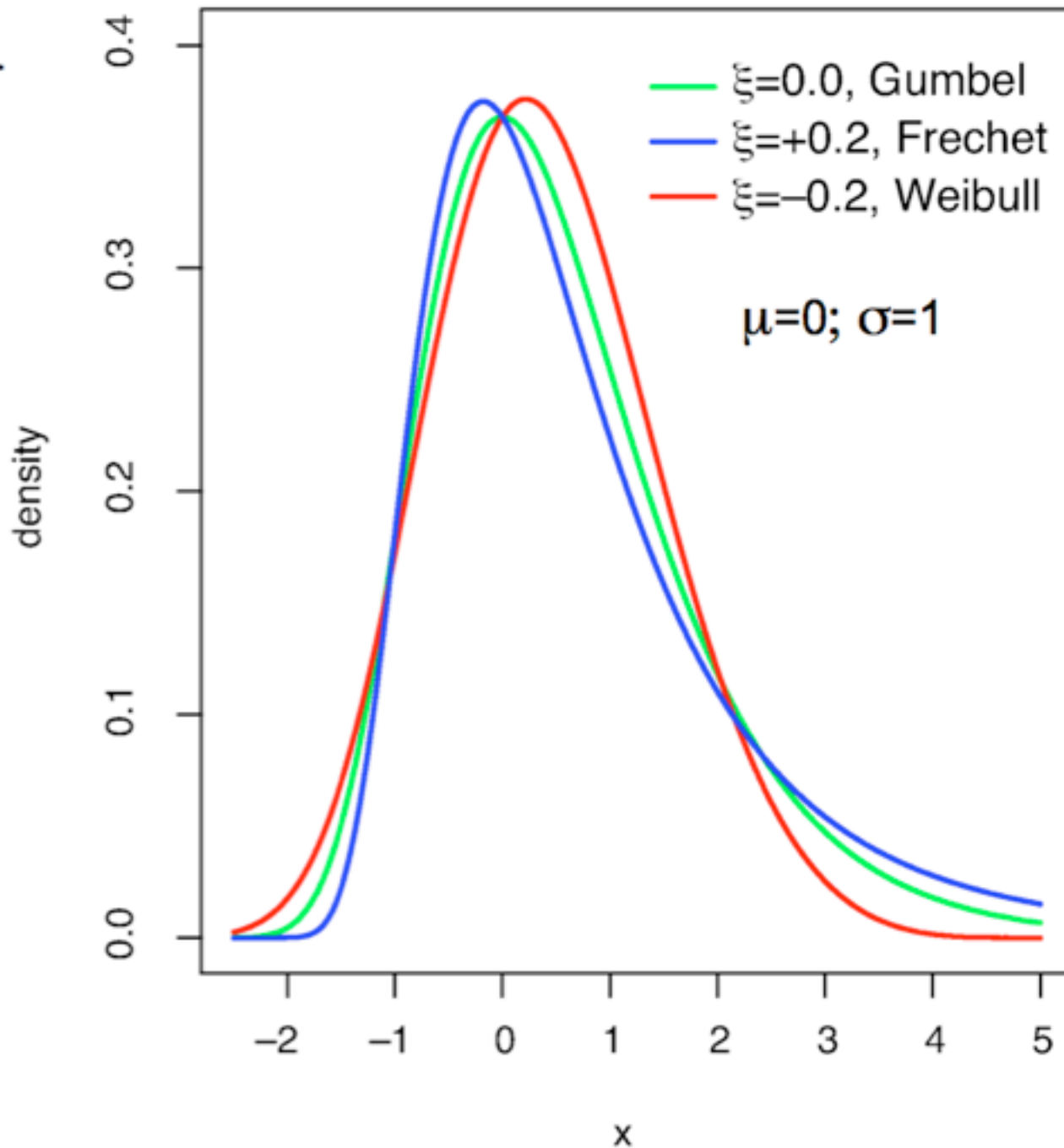
$$G(z) = \exp \left\{ - \left[1 + \xi \left(\frac{z - \mu}{\sigma} \right) \right]^{-1/\xi} \right\}, \quad (2.2)$$

defined on $\{z : 1 + \xi(z - \mu)/\sigma > 0\}$, where $-\infty < \mu < \infty$, $\sigma > 0$ and $-\infty < \xi < \infty$.

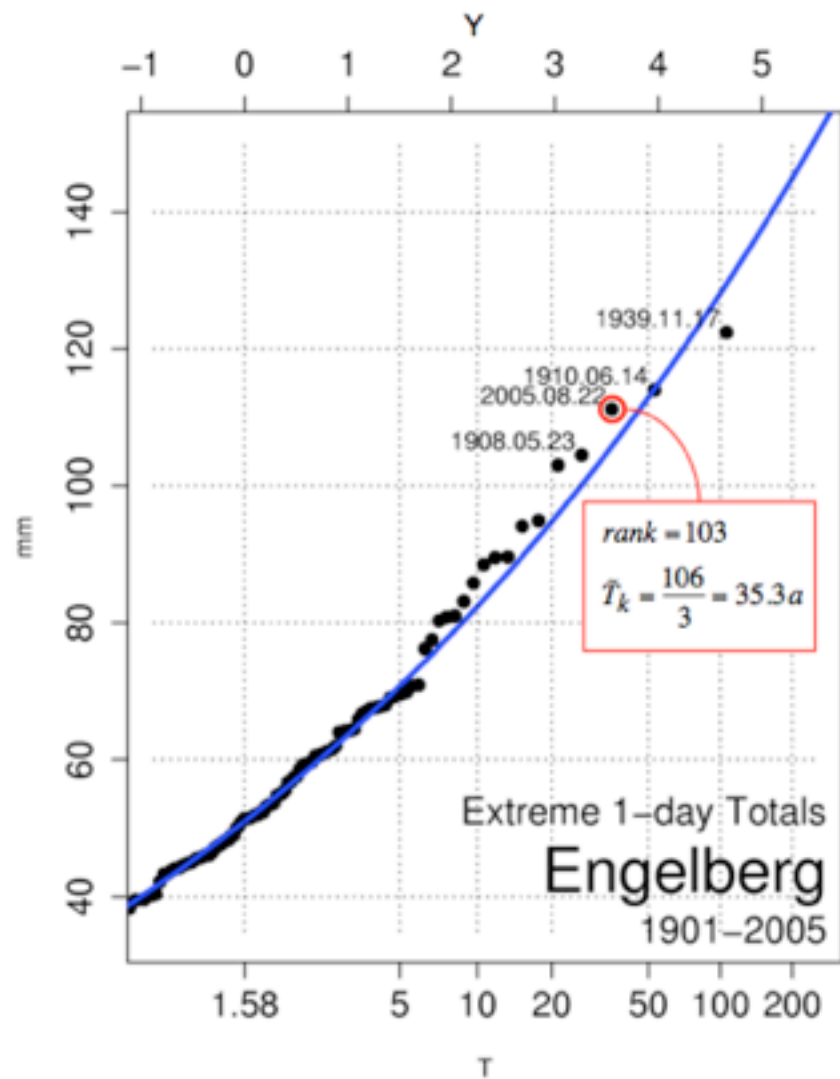
$$z_p = \begin{cases} \mu - \frac{\sigma}{\xi} [1 - \{-\log(1-p)\}^{-\xi}], & \text{for } \xi \neq 0, \\ \mu - \sigma \log \{-\log(1-p)\}, & \text{for } \xi = 0, \end{cases}$$

where $G(z_p) = 1 - p$.

Definition The extreme quantile $z_p = G^{-1}(1 - p)$, where G is the distribution function of M_n , is called the *return level* associated with the *return period* $1/p$.



μ : location
 σ : scale
 ξ : slope



T -axis is transformed such that
Gumbel-Distribution is a straight line.

$F(x) = GEV(x; \mu, \sigma, \xi)$ estimated CDF

$Y(x) = -\log(-\log(F(x)))$ Gumbel Variate

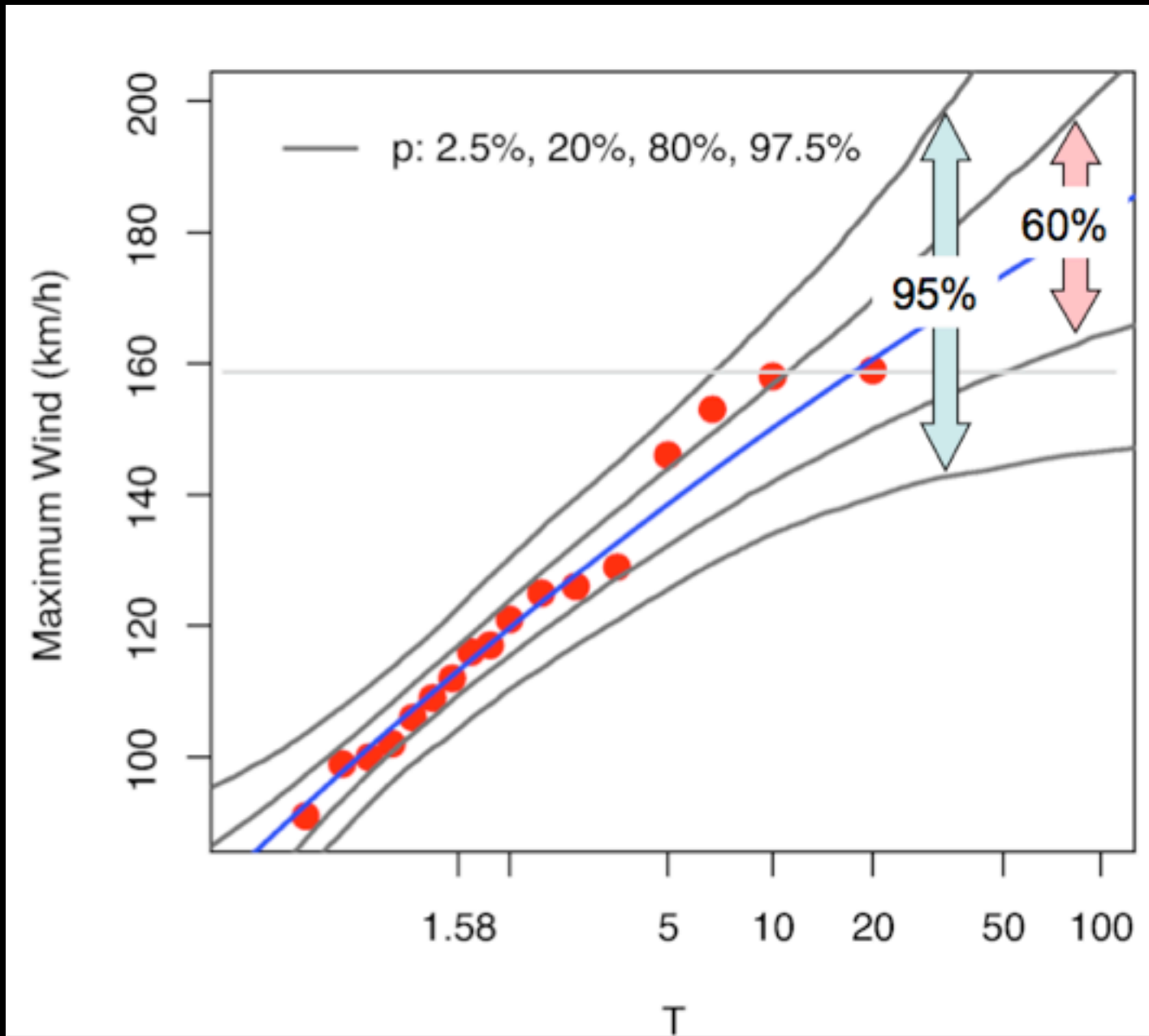
Horizontal axis is linear in Y .

$T(x) = 1/(1 - F(x))$ Return period

x_k $k = 1, \dots, N$ Block Maxima

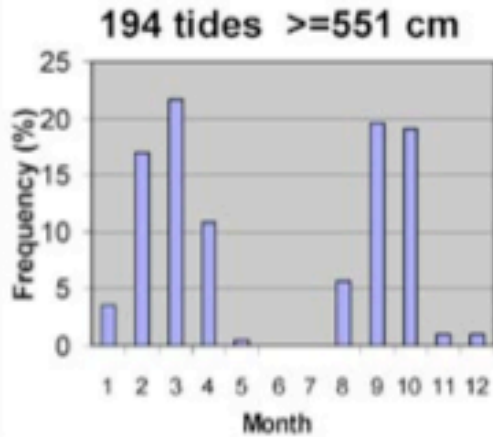
$\tilde{T}_k = \frac{N+1}{N+1 - rank(x_k)}$ plotting points of
block maxima x_k

Parametric resampling and confidence intervals

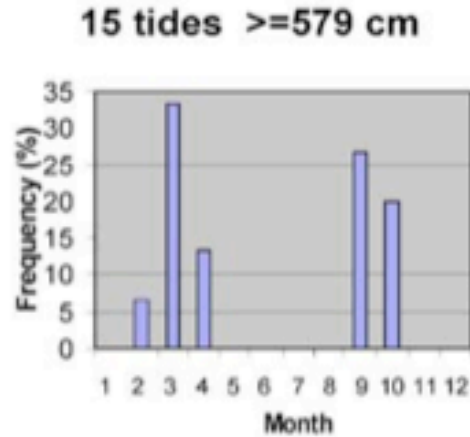


Extreme tides: temporal series

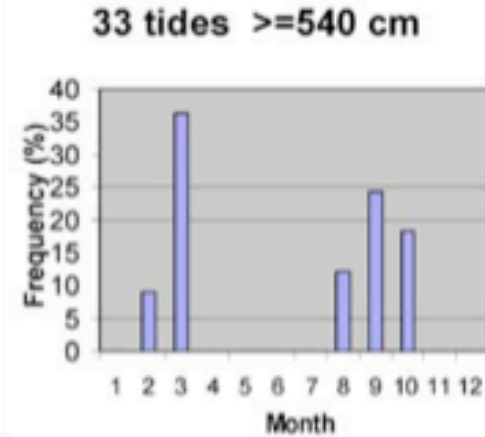
9. Port-Tudy



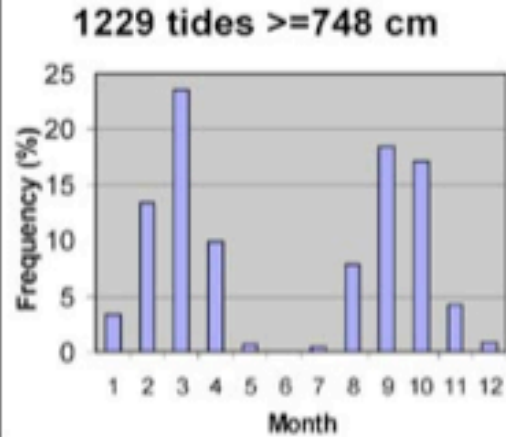
10. Le Crouesty



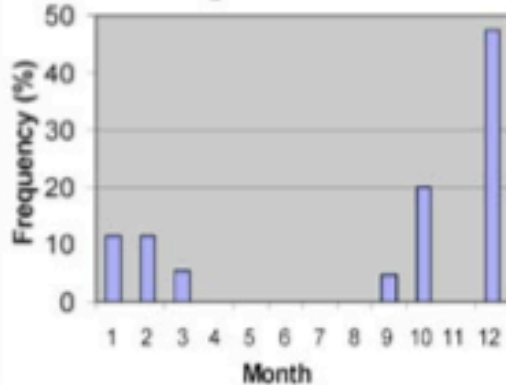
11. Concarneau



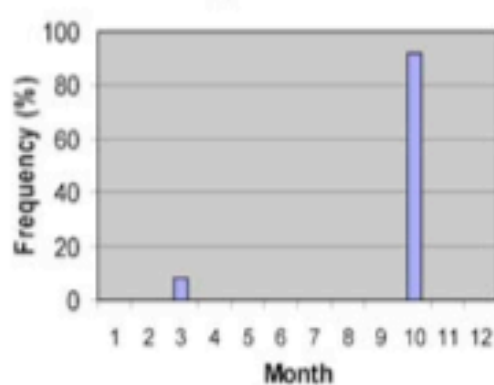
12. Brest



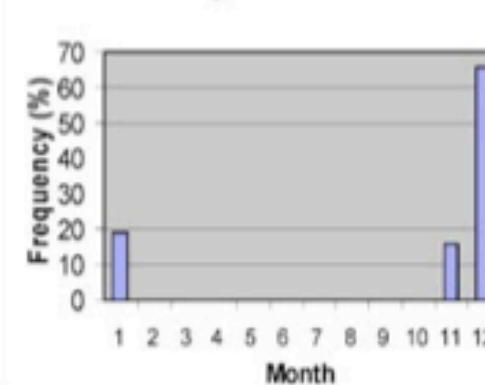
193 surges ≥ 62 cm



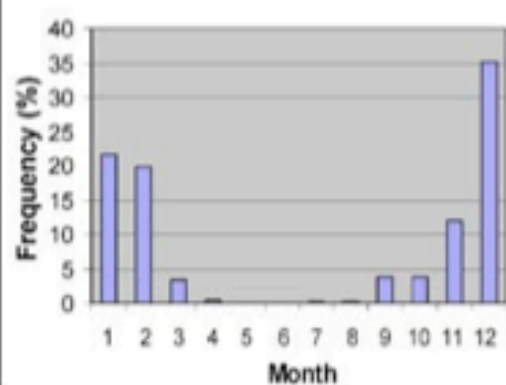
12 surges ≥ 51 cm



32 surges ≥ 72 cm



1186 surges ≥ 63 cm



Monthly matching possibilities: 45.08%

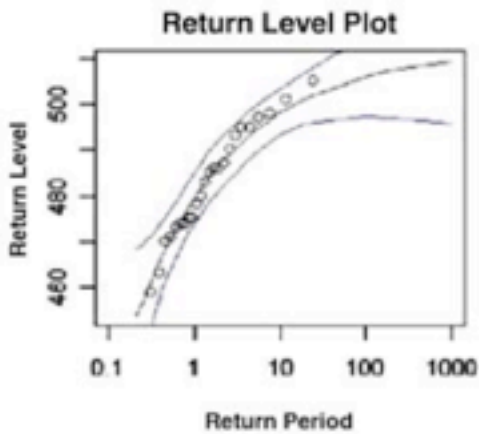
Monthly matching possibilities: 28.33%

Monthly matching possibilities: 0%

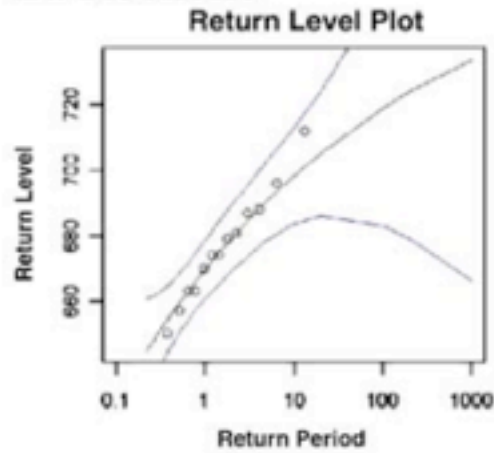
Monthly matching possibilities: 30,26%

Expected return periods

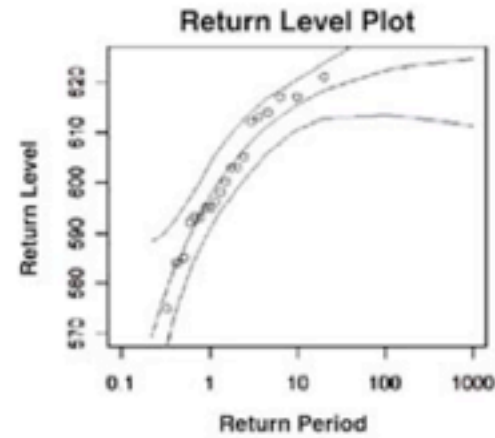
Boucau



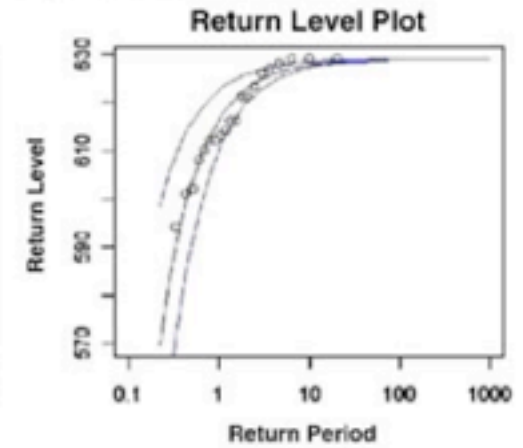
La Rochelle



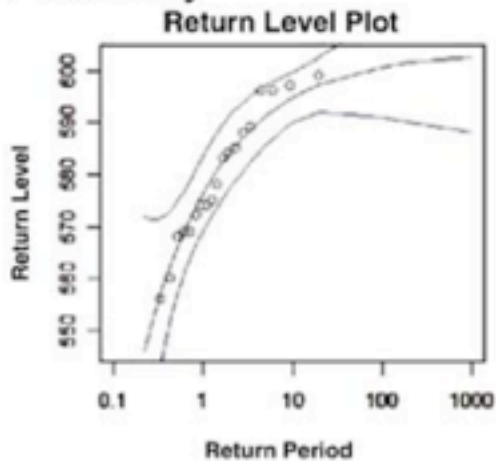
Saint-Gildas



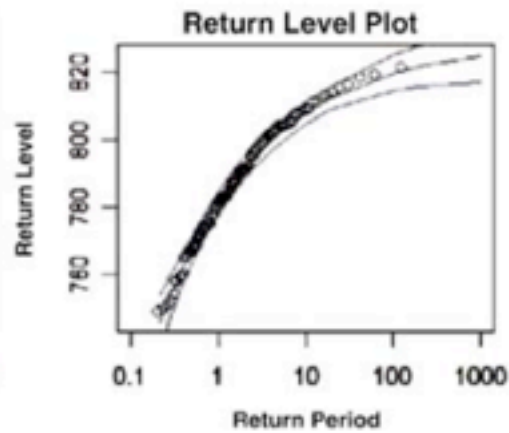
Saint-Nazaire



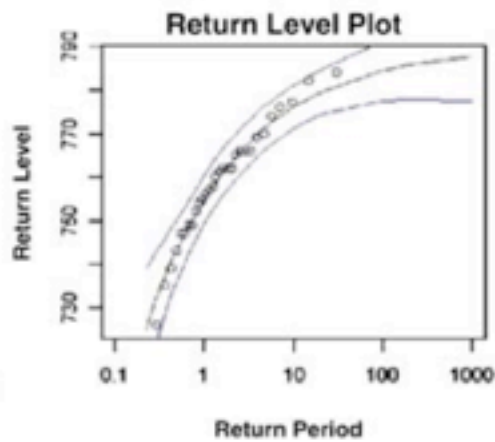
Port Tudy



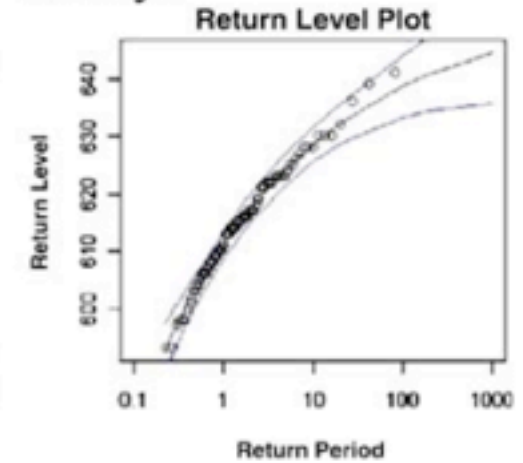
Brest



Le Conquet



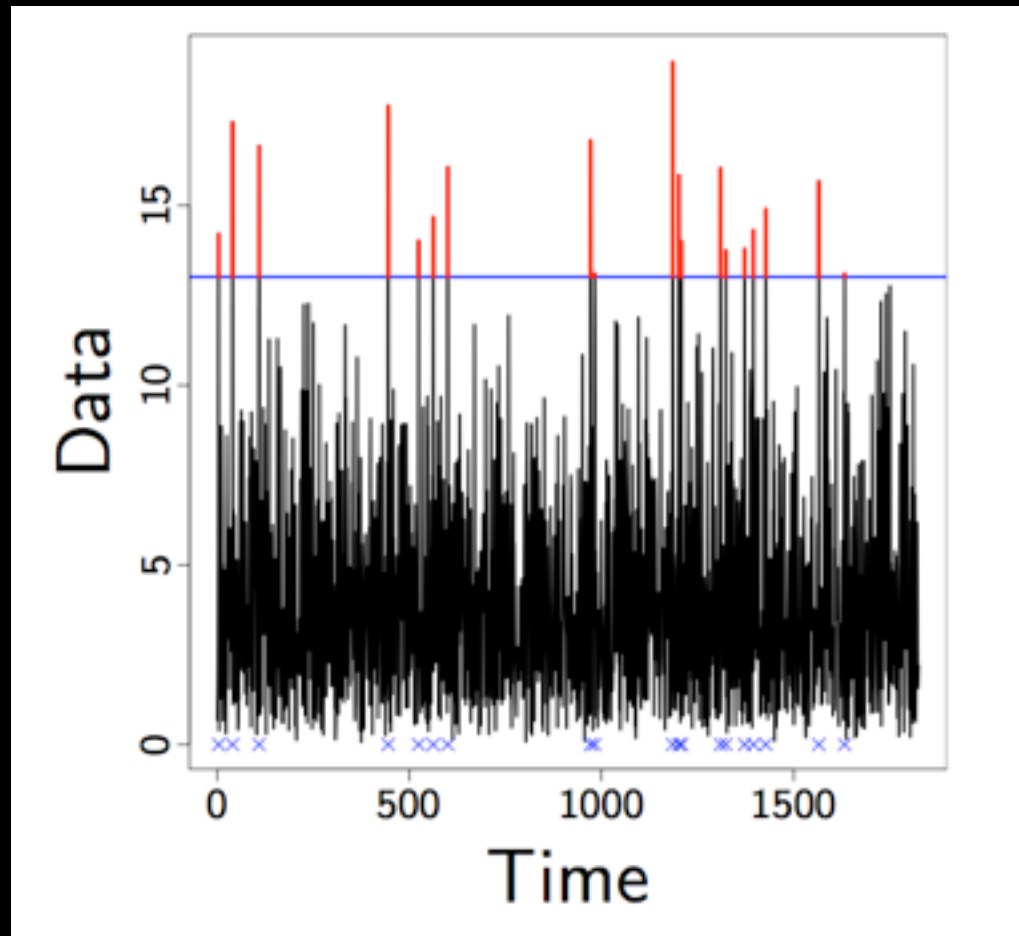
Newlyn



Basics of extreme value theory IV

Estimation of GEV parameters:

- ★ Method of moments: not robust and VERY unstable
- ★ **Block maxima**: extreme events in FIXED intervals
- ★ **Peak over thresholds**: all extreme events ABOVE a fixed value



Basics of extreme value theory V

Conditional excess distribution function for a threshold level u :

$$F_u(y) = P(y \geq X - u | X > u), \quad 0 \leq y < \infty$$

$$F_u(y) = \frac{F(u + y) - F(u)}{1 - F(u)}, \quad y > 0.$$

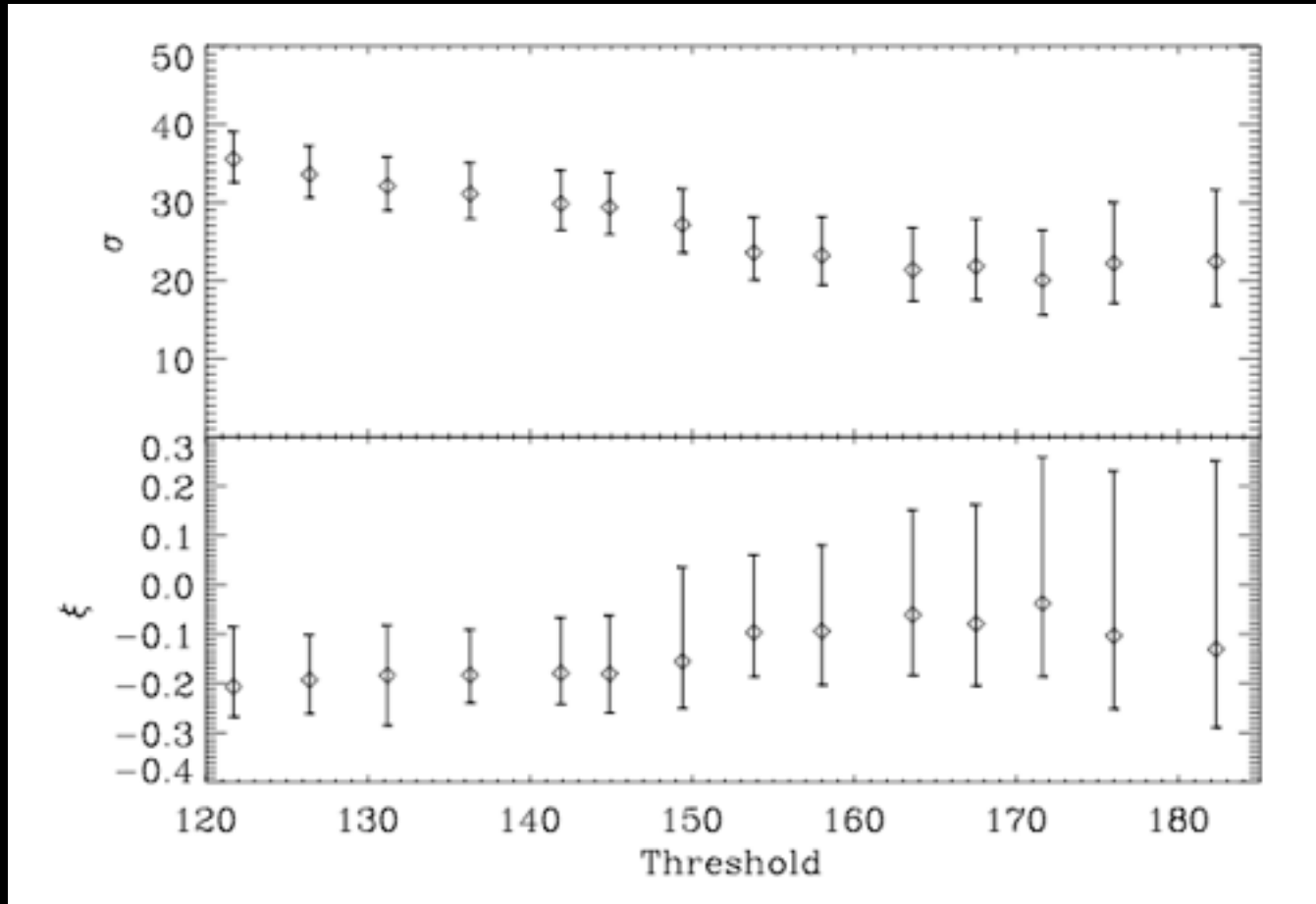
Theorem (Pickands 1975):

$$F_u(y) \approx \begin{cases} 1 - \left(1 + \frac{\xi}{\sigma}y\right)^{-1/\xi} & \text{if } \xi \neq 0 \\ 1 - e^{-y/\sigma} & \text{if } \xi = 0, \end{cases}$$

so that the cumulative DF for events above u , taking $x=u+y$, is :

$$F(x) = 1 - \frac{N_u}{n} \left[1 + \frac{\xi}{\sigma}(x - u)\right]^{-1/\xi}$$

For a given sample, calculate the likelihood and apply standard Bayesian estimation :



Star Formation and the Origins of the Stellar Initial Mass Function



Why are there fewer massive stars than low-mass stars?



- Typical stellar mass \sim Jeans Mass.

$$M_{\text{Jeans}} \equiv \left[\frac{5 R_g T}{2 G \mu} \right]^{\frac{3}{2}} \left[\frac{4\pi}{3} \rho \right]^{-\frac{1}{2}} .$$

- In Molecular clouds:
 - Temperatures ~ 10 K,
 - Lots of structure,
 - Dense cores: $\rho \sim 10^{-19}$ g/cm³.
- \Rightarrow masses $M \sim 0.7 M_{\text{Sol}}$.
 - Agrees well with observations:
 - $M_{\text{median}} \sim 0.5 M_{\text{Sol}}$.
- Not all stars have the same mass \Rightarrow Distribution?
- **The IMF!**

- **Salpeter-IMF** (power-law): $\xi(m) = k m^{-2.35}$, originally only for 0.4 to 10 M_{Sol} .

- **Miller-Scalo IMF** (log-normal):

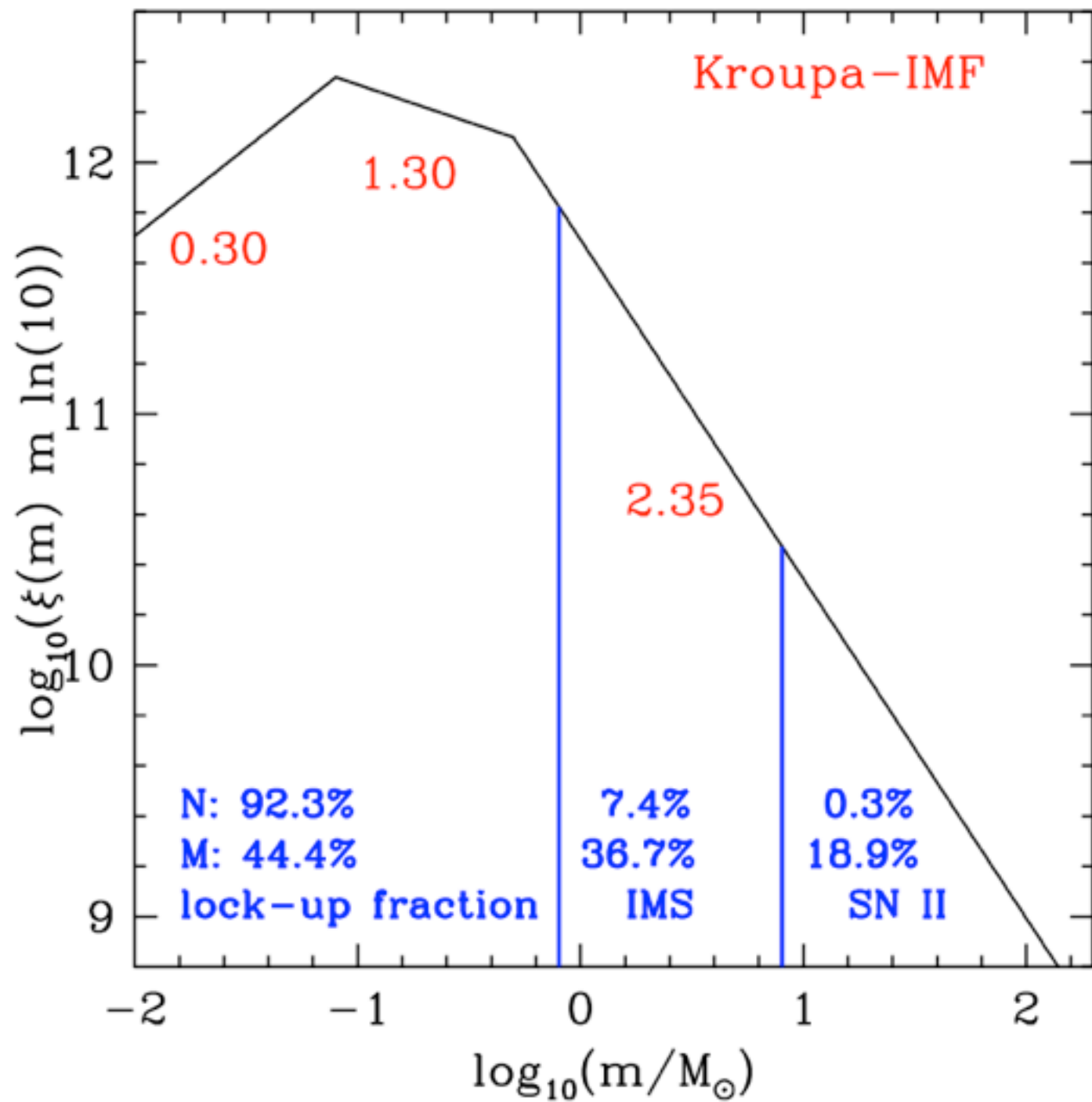
$$\xi(lm) = k e^{-\frac{(lm+1.02)^2}{0.9248}},$$

with $lm = \log_{10} m$.

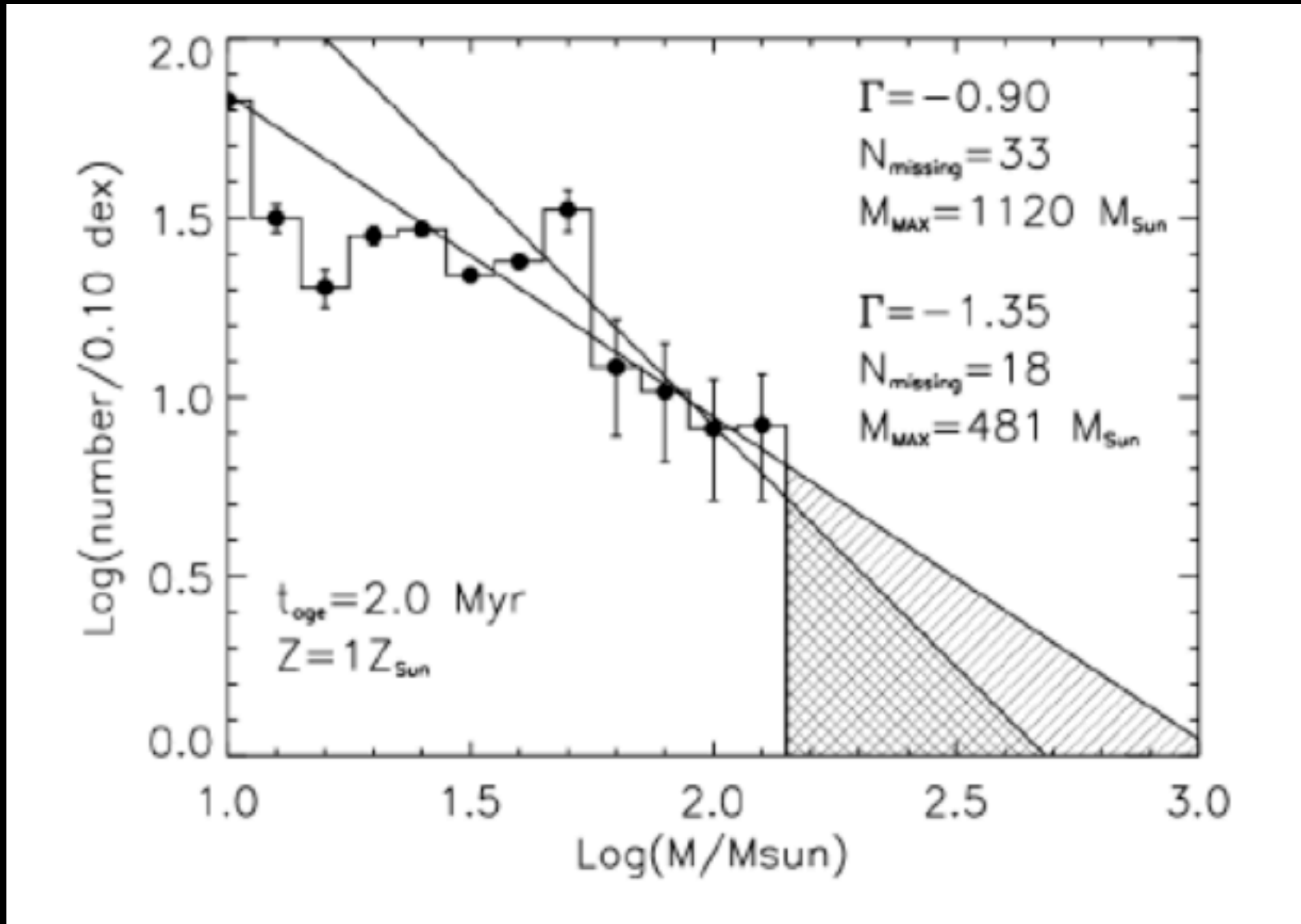
- **Kroupa-IMF** (multi power-law):

$$\xi(m) = k \begin{cases} \left(\frac{m}{m_H}\right)^{-\alpha_0} & , m_{\text{low}} \leq m < m_H, \\ \left(\frac{m}{m_H}\right)^{-\alpha_1} & , m_H \leq m < m_0, \\ \left(\frac{m_0}{m_H}\right)^{-\alpha_1} \left(\frac{m}{m_0}\right)^{-\alpha_2} & , m_0 \leq m < m_1, \\ \left(\frac{m_0}{m_H}\right)^{-\alpha_1} \left(\frac{m_1}{m_0}\right)^{-\alpha_2} \left(\frac{m}{m_1}\right)^{-\alpha_3} & , m_1 \leq m < m_{\text{max}}, \end{cases}$$

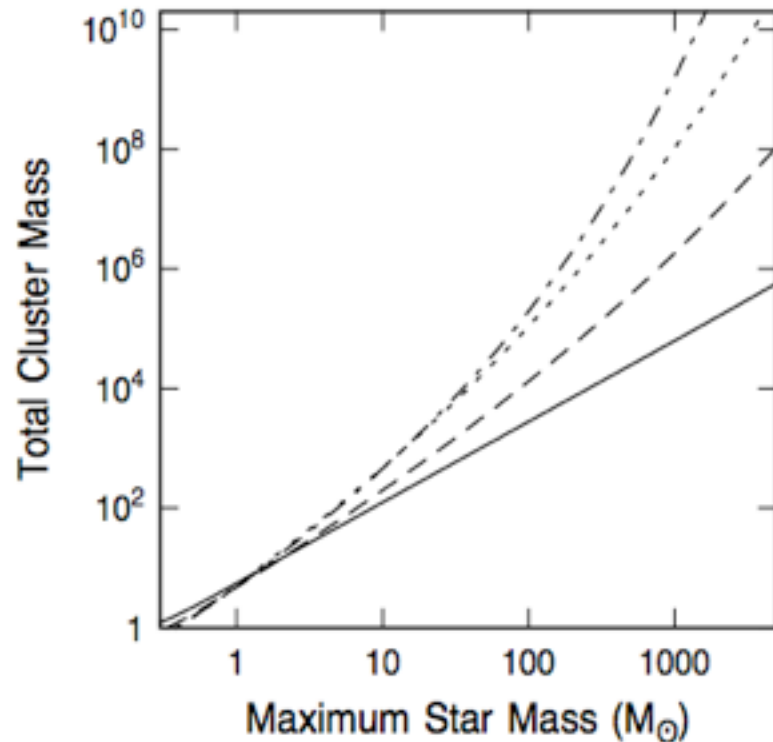
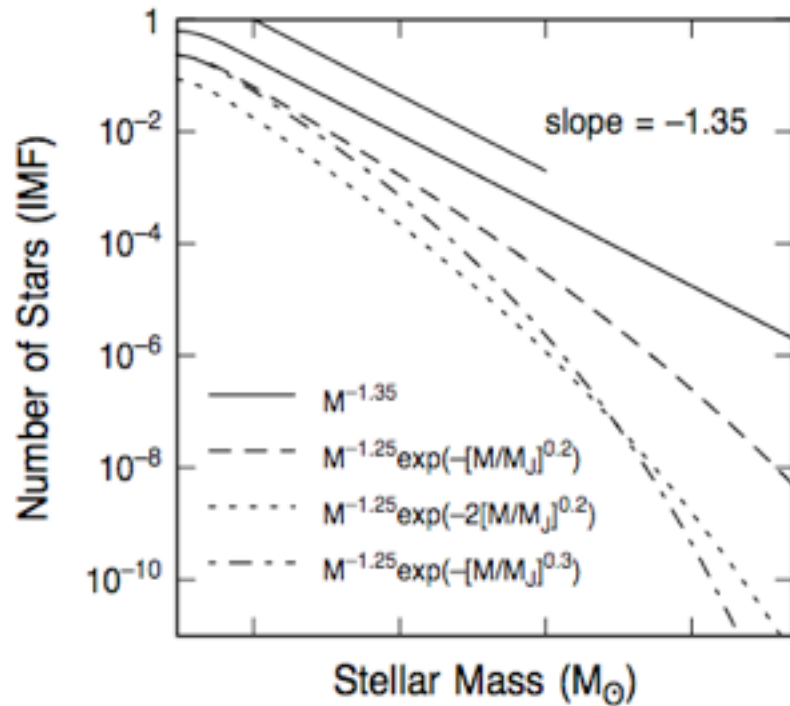
$$\begin{aligned} \alpha_0 &= +0.30 & , & 0.01 \leq m/M_{\odot} < 0.08, \\ \alpha_1 &= +1.30 & , & 0.08 \leq m/M_{\odot} < 0.50, \\ \alpha_2 &= +2.35 & , & 0.50 \leq m/M_{\odot} < 1.00, \\ \alpha_3 &= +2.35 & , & 1.00 \leq m/M_{\odot}. \end{aligned}$$



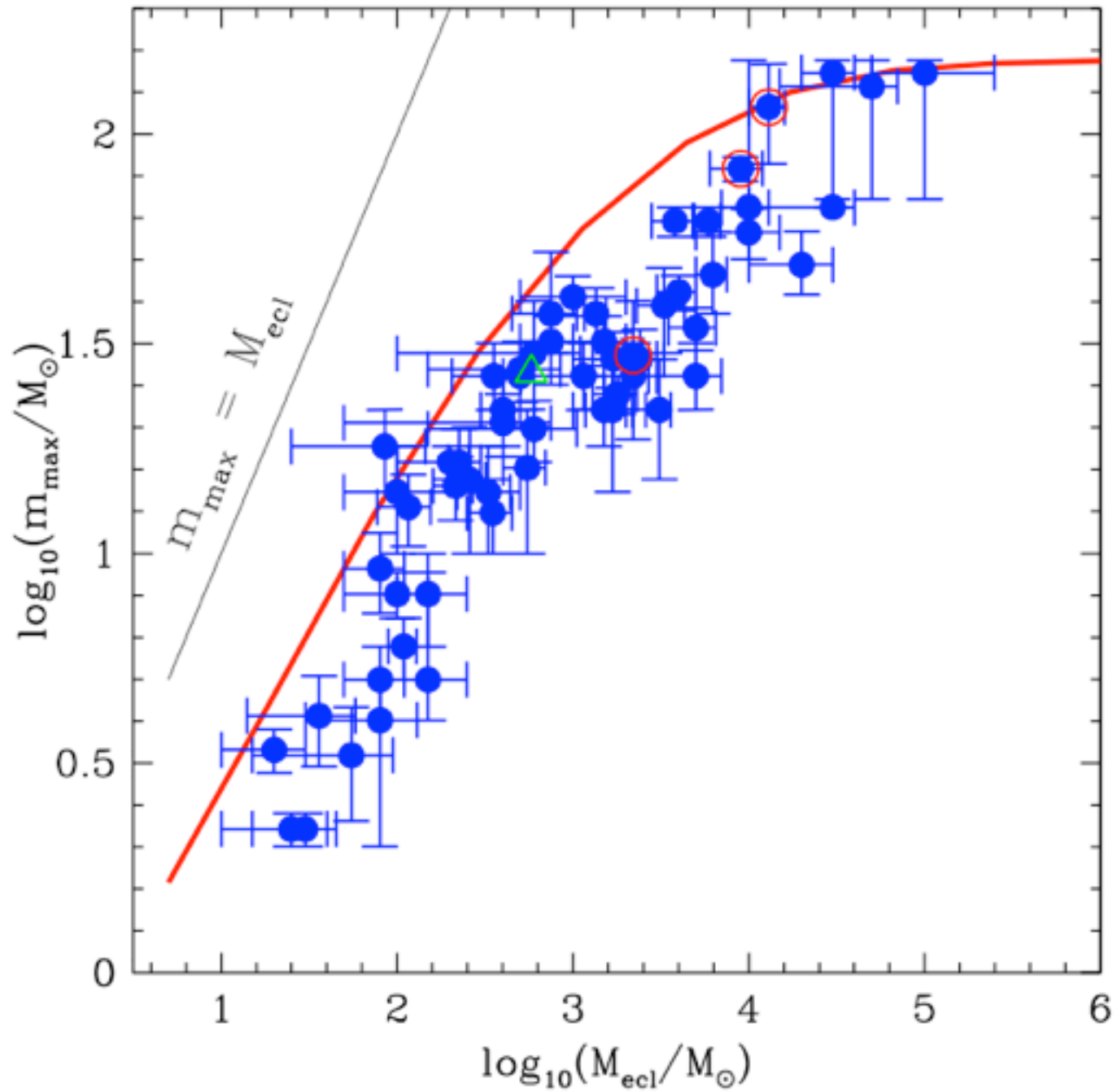
Evidence for an upper mass cutoff around $150 M_{\odot}$?



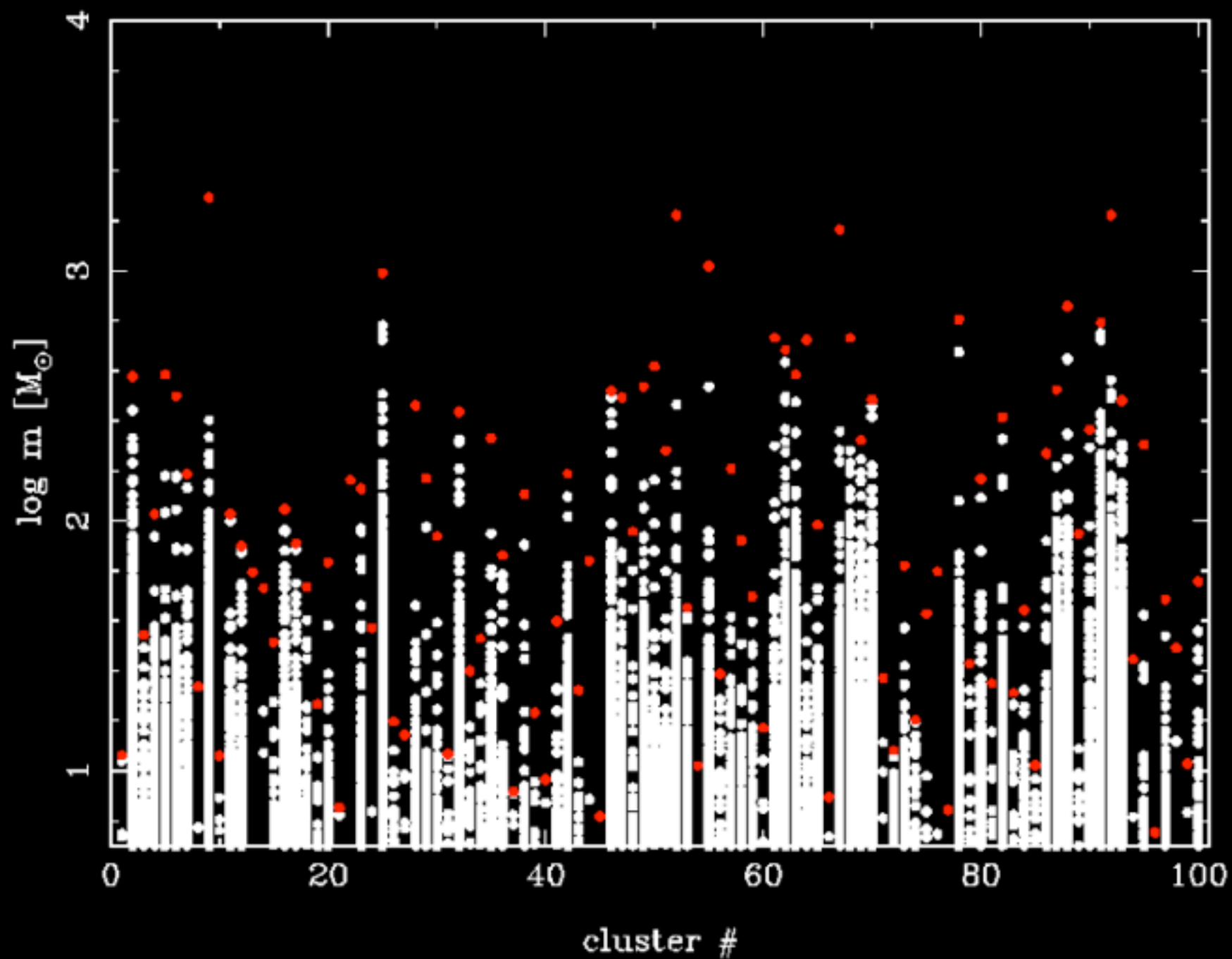
Figier (2005) Nature 434, 592



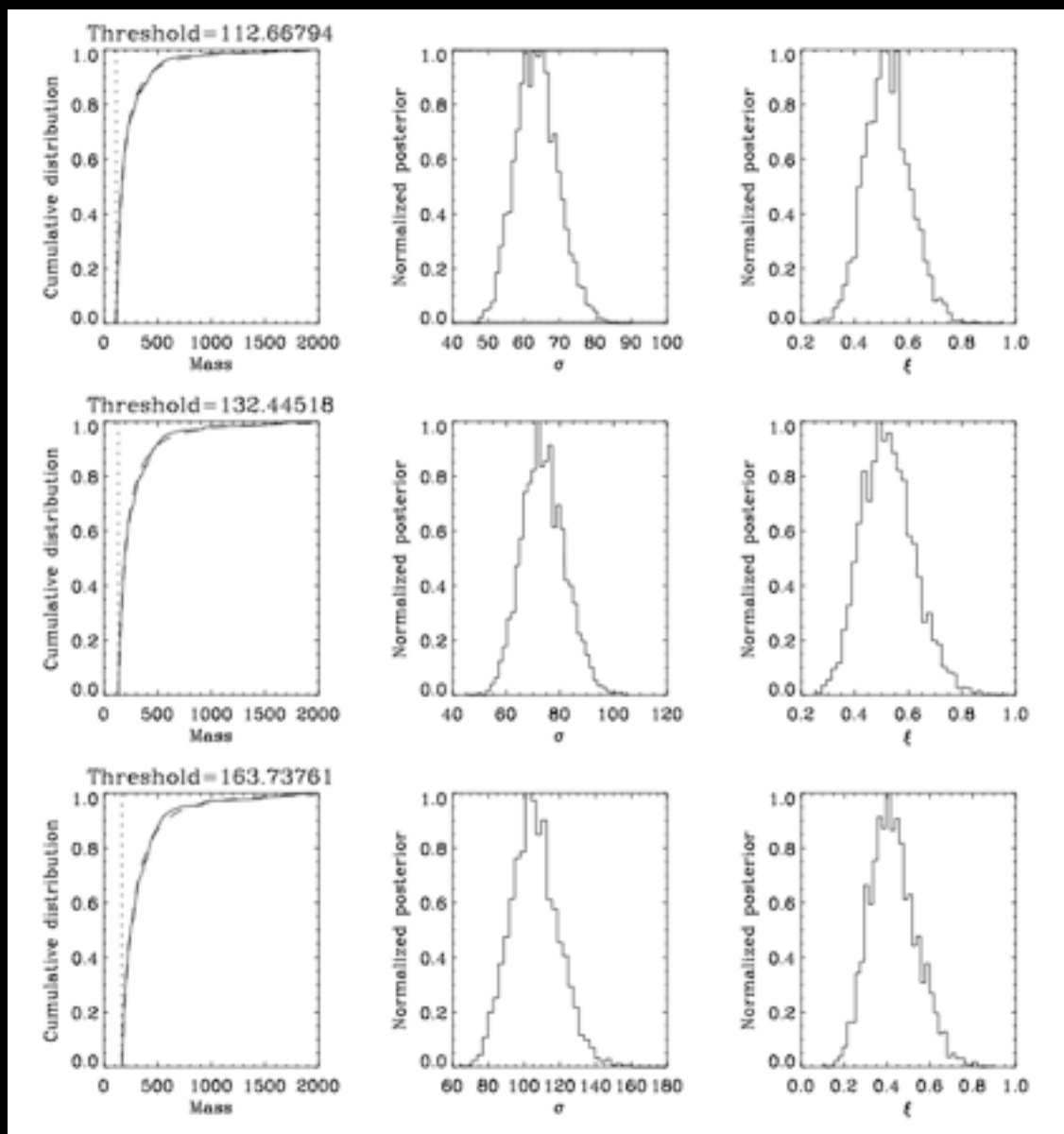
One expects a correlation
between the total
cluster mass
and the maximum
observed stellar mass



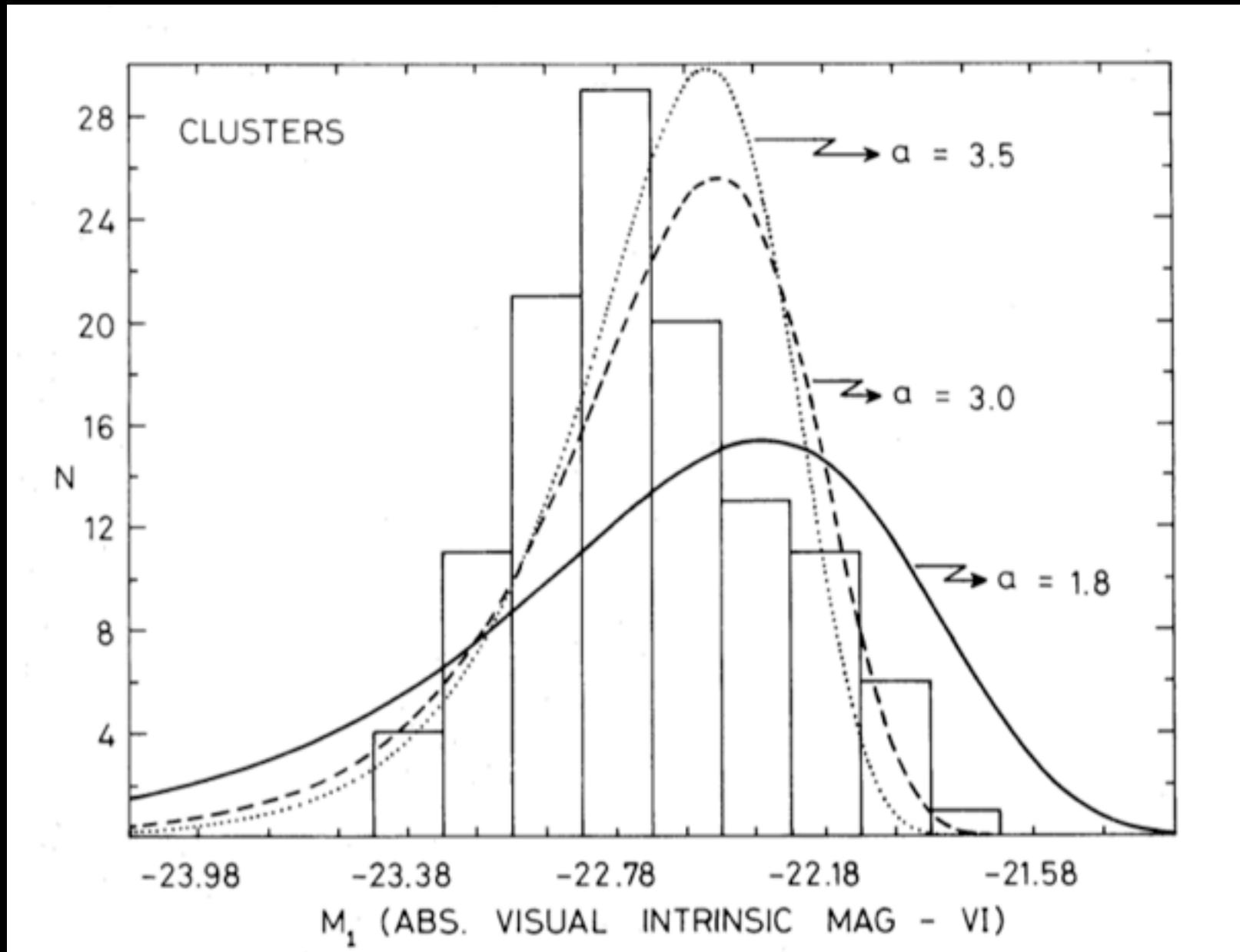
sampled_IMF_a03.out



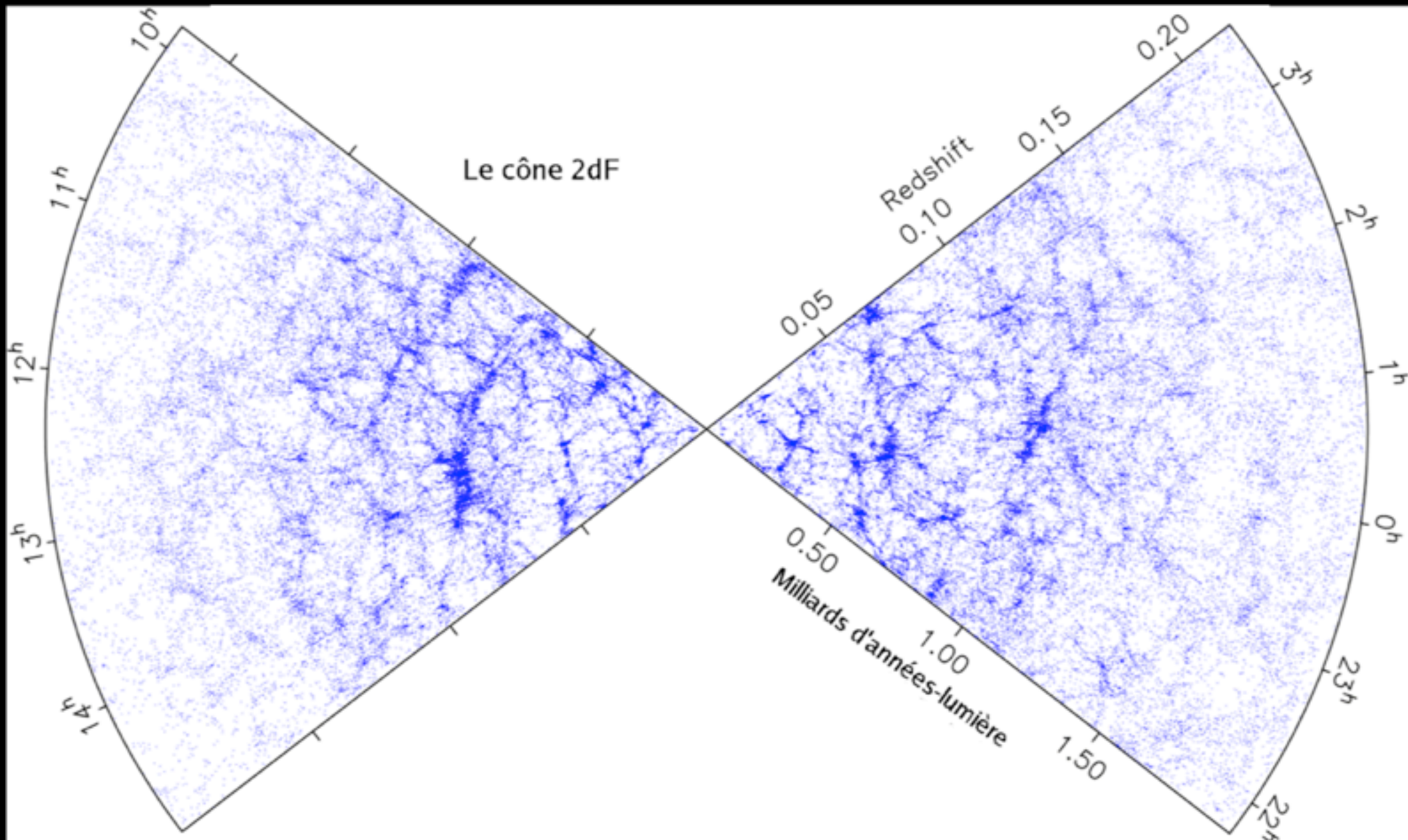
Posterior PDFs for scale and slope of IMF for massive stars



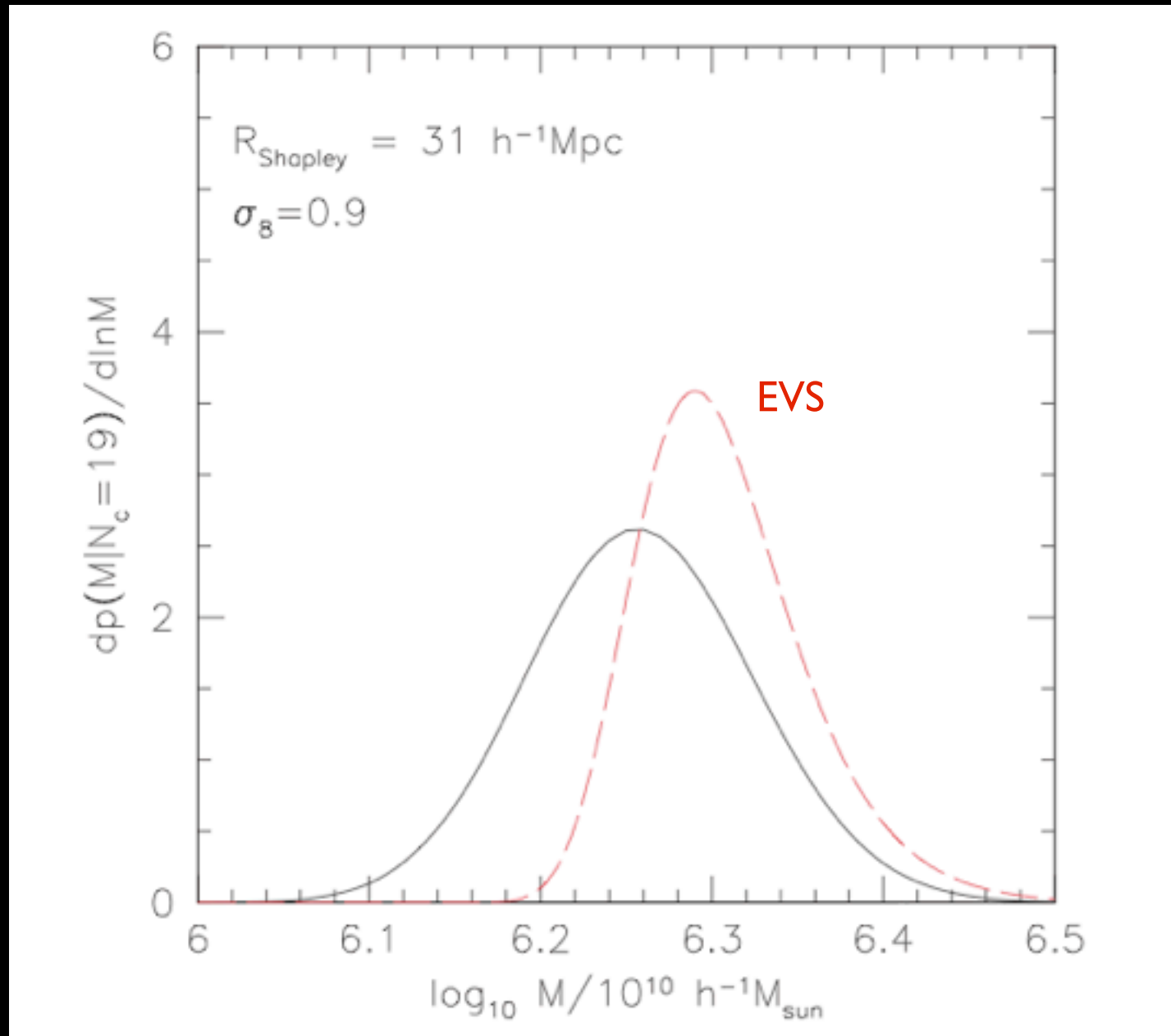
Distribution function of brightest galaxies in clusters



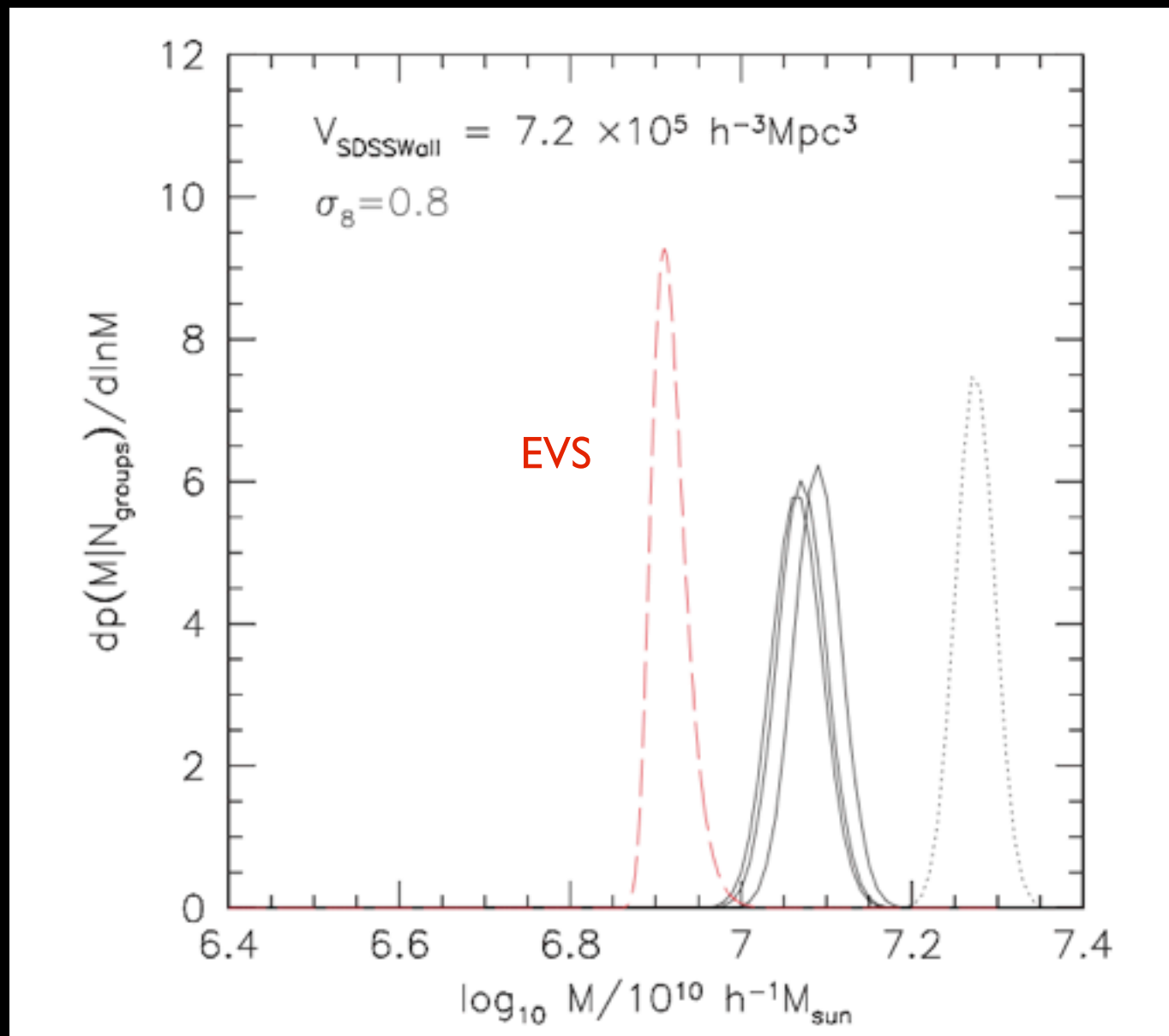
Probability of finding the largest structures in a cosmological volume

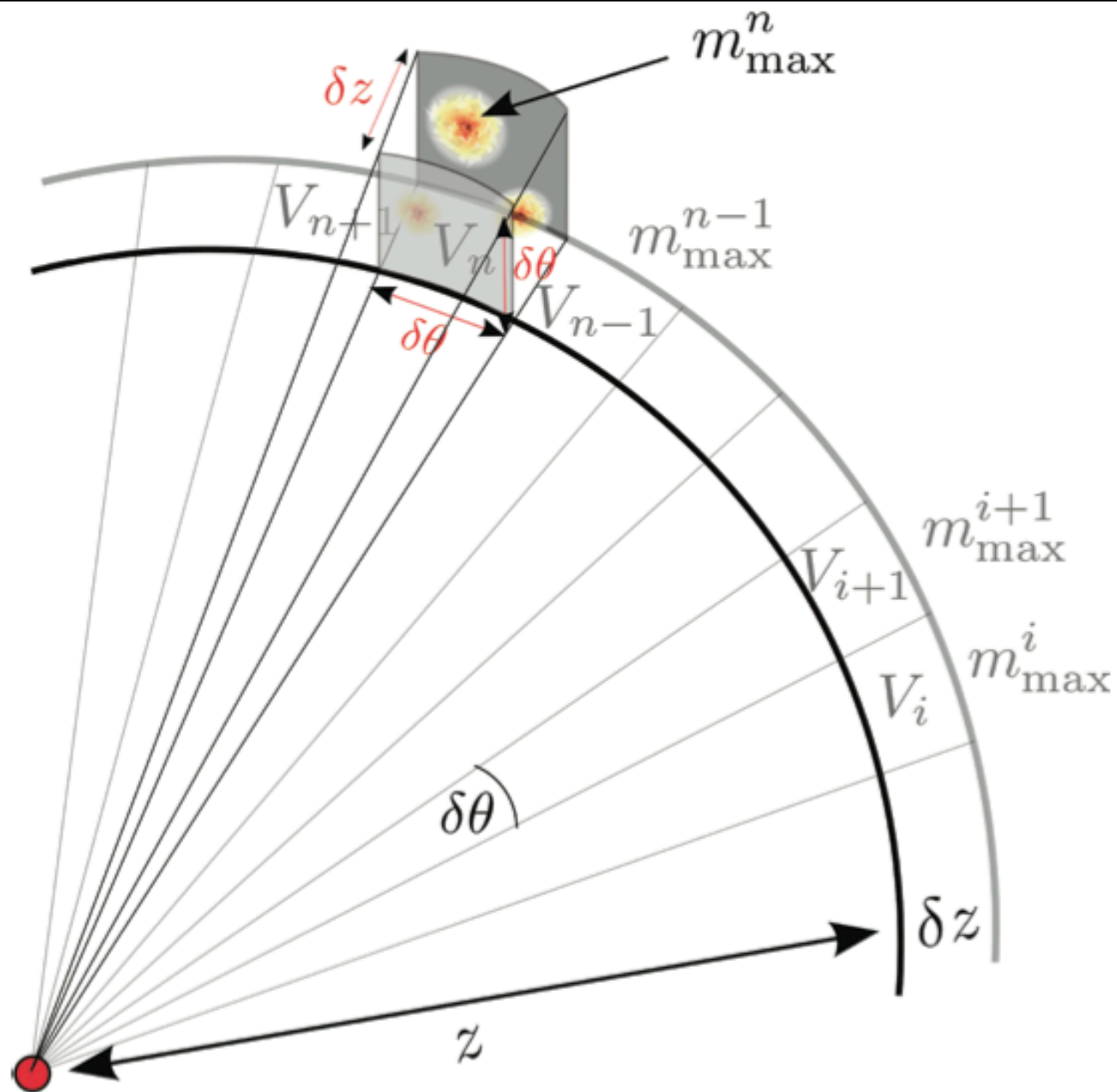


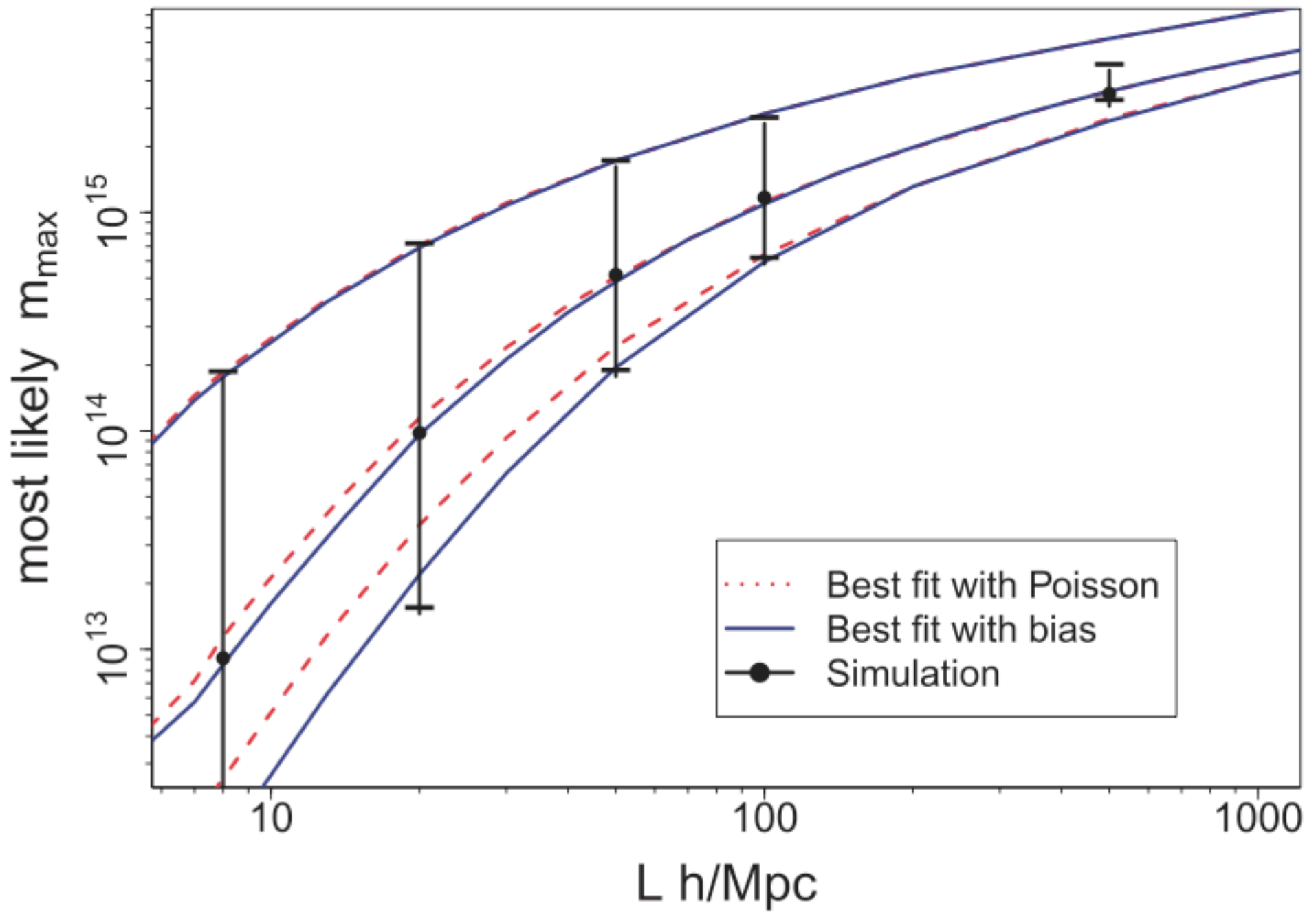
Probability of finding a Shapley supercluster



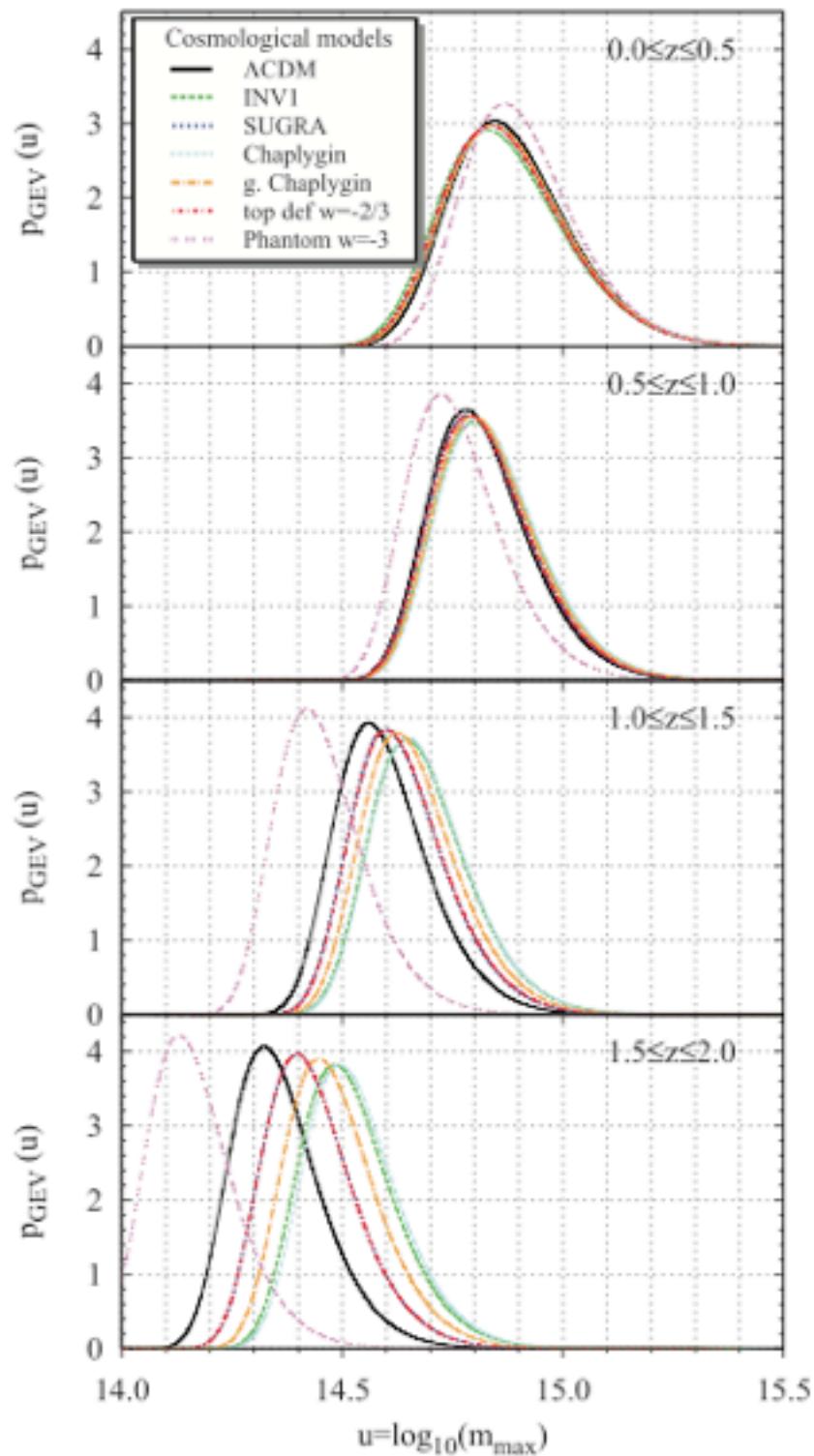
Probability of finding the SDSS Great Wall



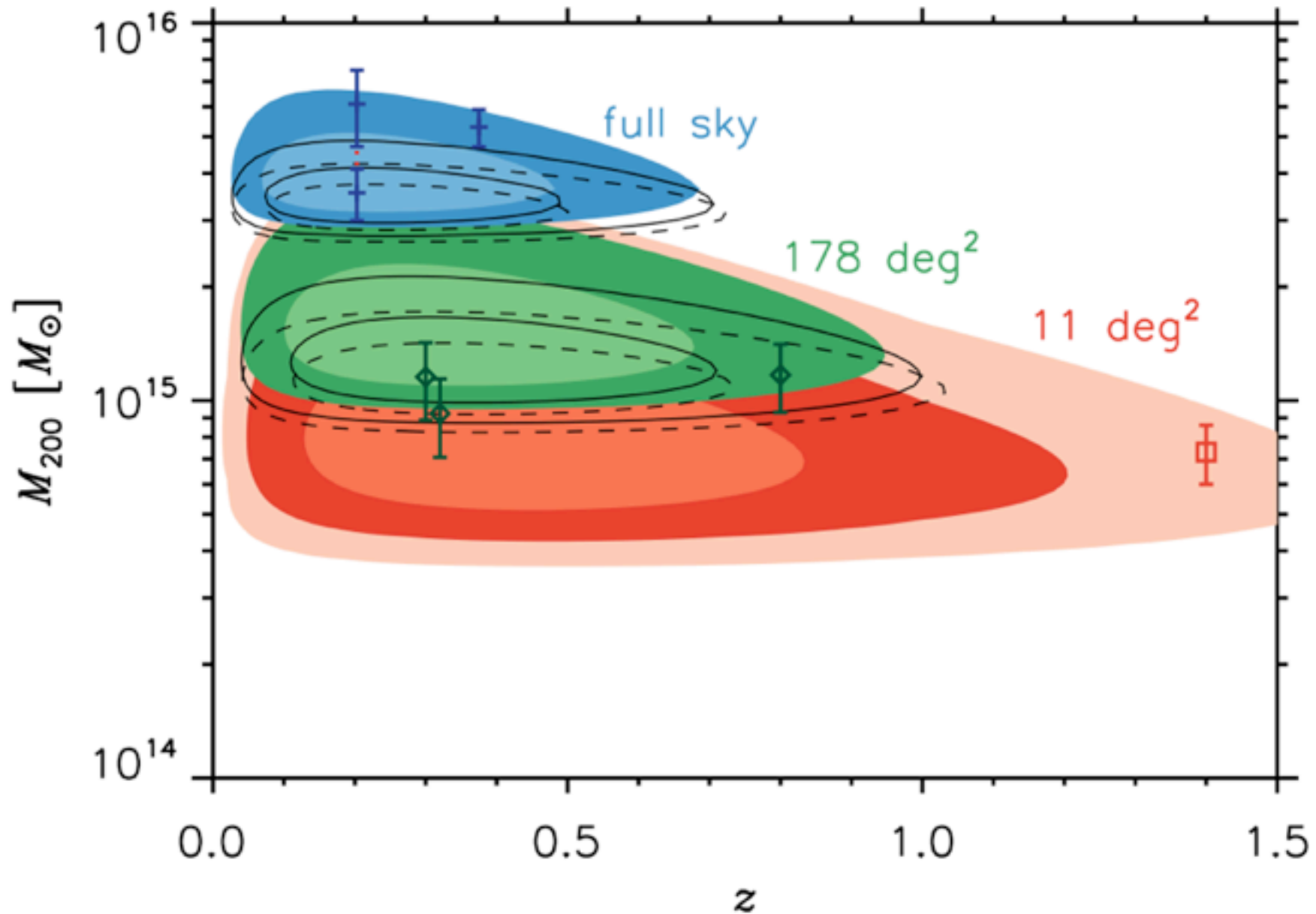




PDF of the most massive clusters of galaxies as a function of redshift and equation of state of Dark Energy

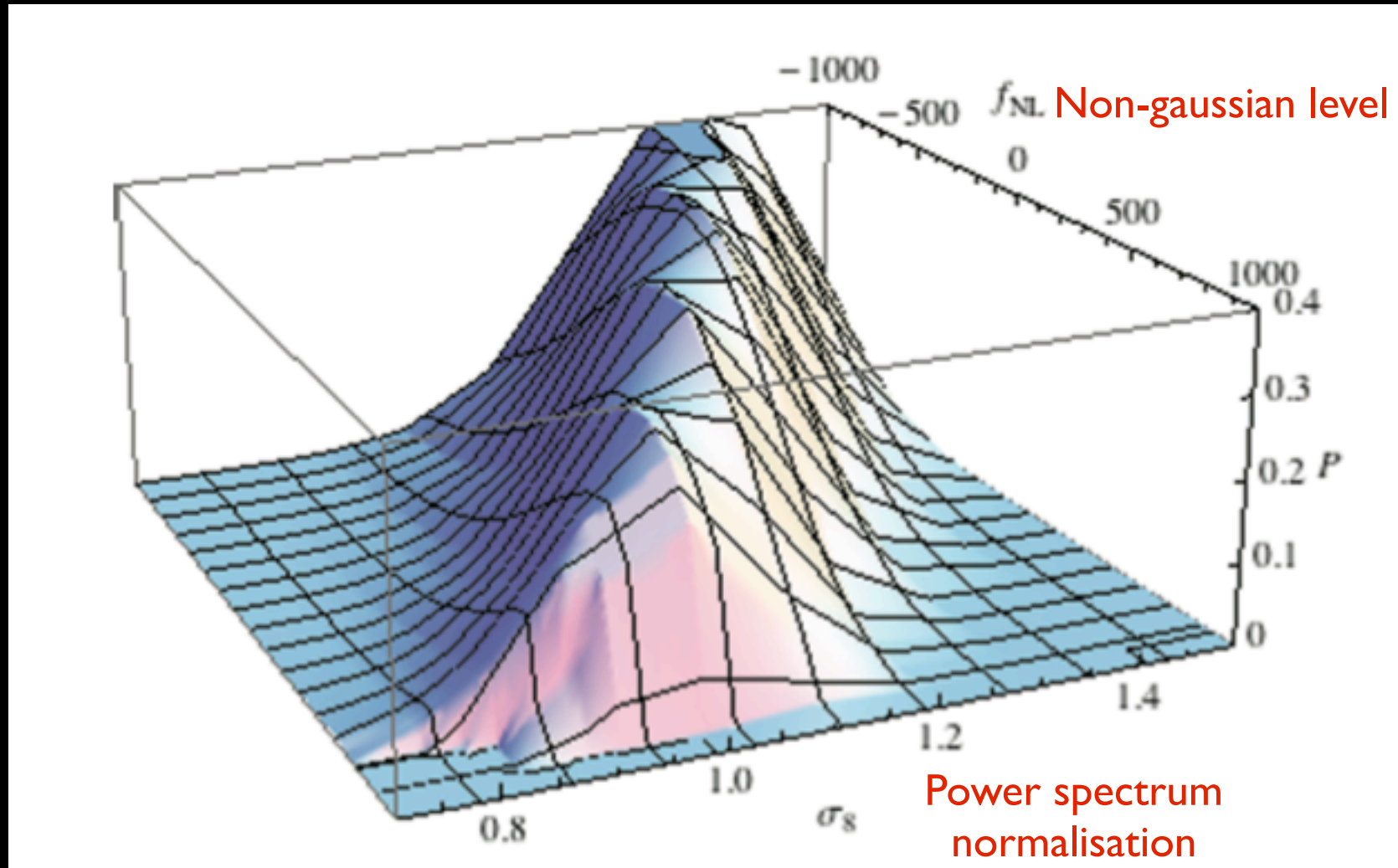


Most massive galaxy clusters expected in a given survey

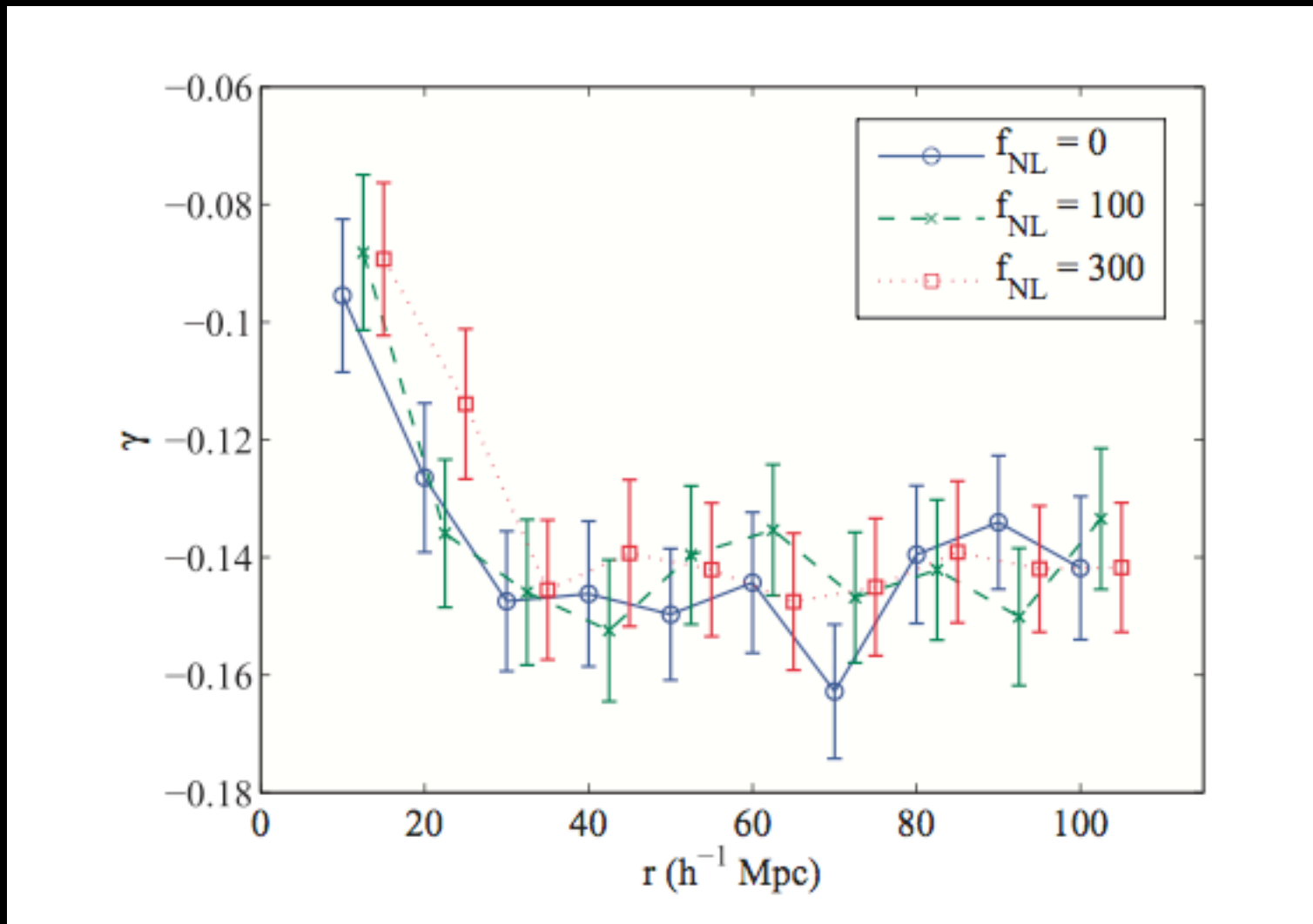


XMMU J2235.3-2557: a massive cluster

$$M_{324} = (6.4 \pm 1.2) 10^{14} M_{\odot} \text{ at } z=1.4$$



.... but the shape parameter of the EVS of the halo mass function does not discriminate non-gaussianity



Testing General Relativity

Abundance

$$n(M, z) = \int_0^M f(\sigma) \frac{\bar{\rho}_m}{M'} \frac{d \ln \sigma^{-1}}{dM'} dM'$$

$$\sigma^2(M, z) = \frac{1}{2\pi^2} \int_0^\infty k^2 P(k, z) |W_M(k)|^2 dk$$

$$P(k, z) \propto k^{n_s} T^2(k, z_t) D(z)^2$$

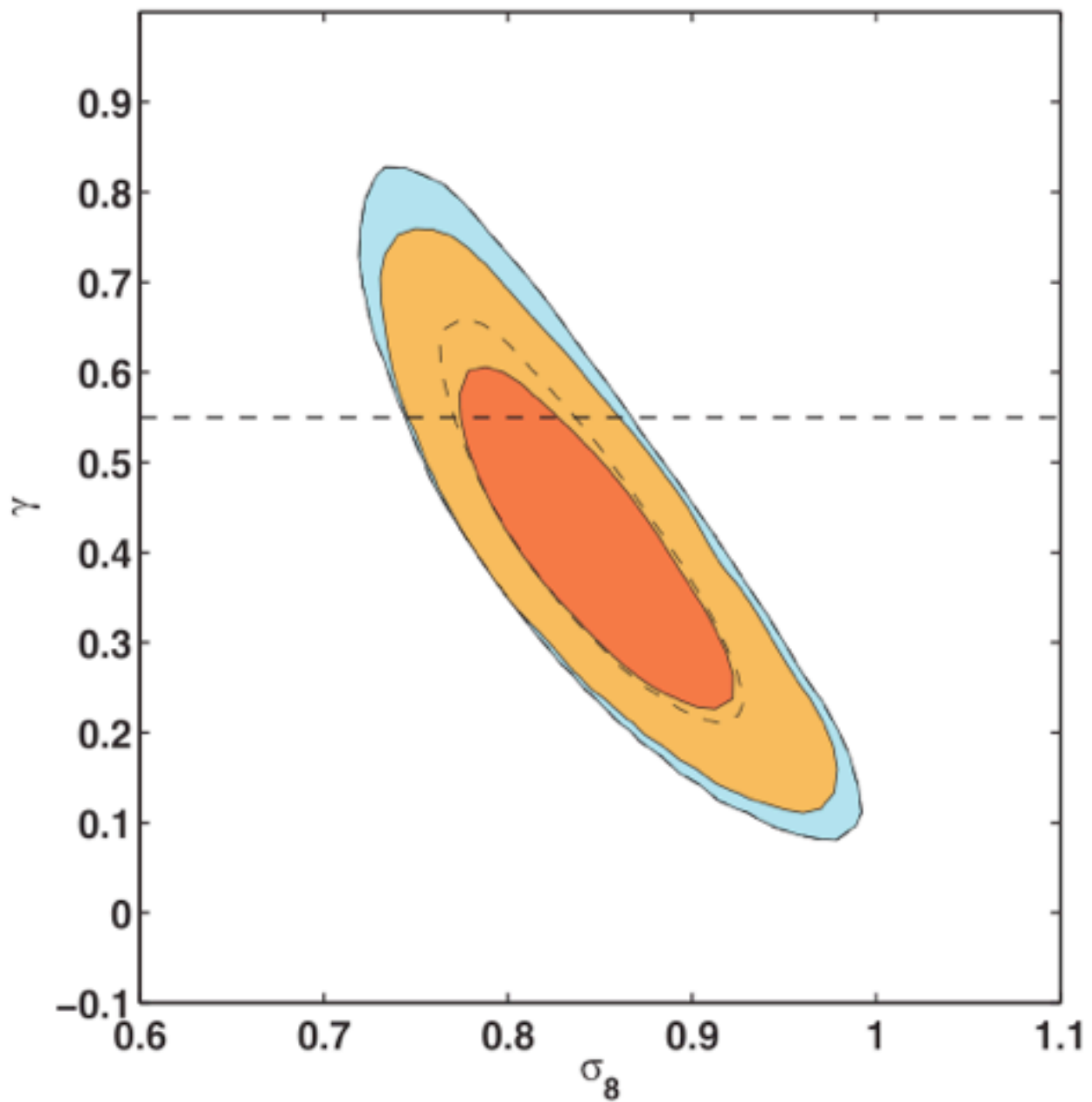
$$D(z) \equiv \frac{\delta(z)}{\delta(z_t)}$$

$$\frac{d \ln \delta}{d \ln a} = \Omega_m(a)^\gamma$$

$$\Omega_m(a) = \Omega_m a^{-3} / E(a)^2$$

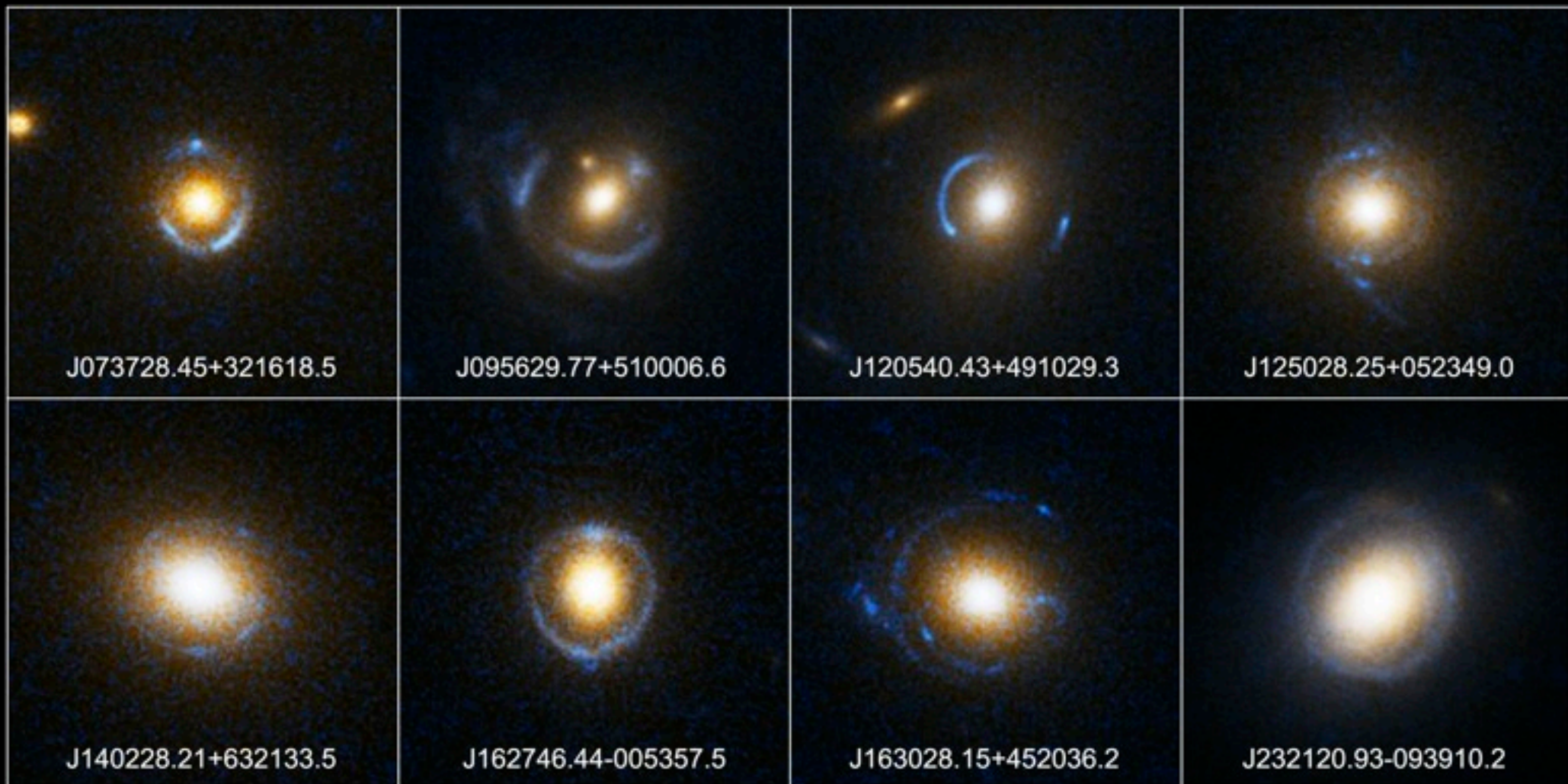
GR $\gamma \sim 0.55$

$$E(a) = \left[\Omega_m a^{-3} + \Omega_{de} a^{-3(1+w)} + \Omega_k a^{-2} \right]^{1/2}$$



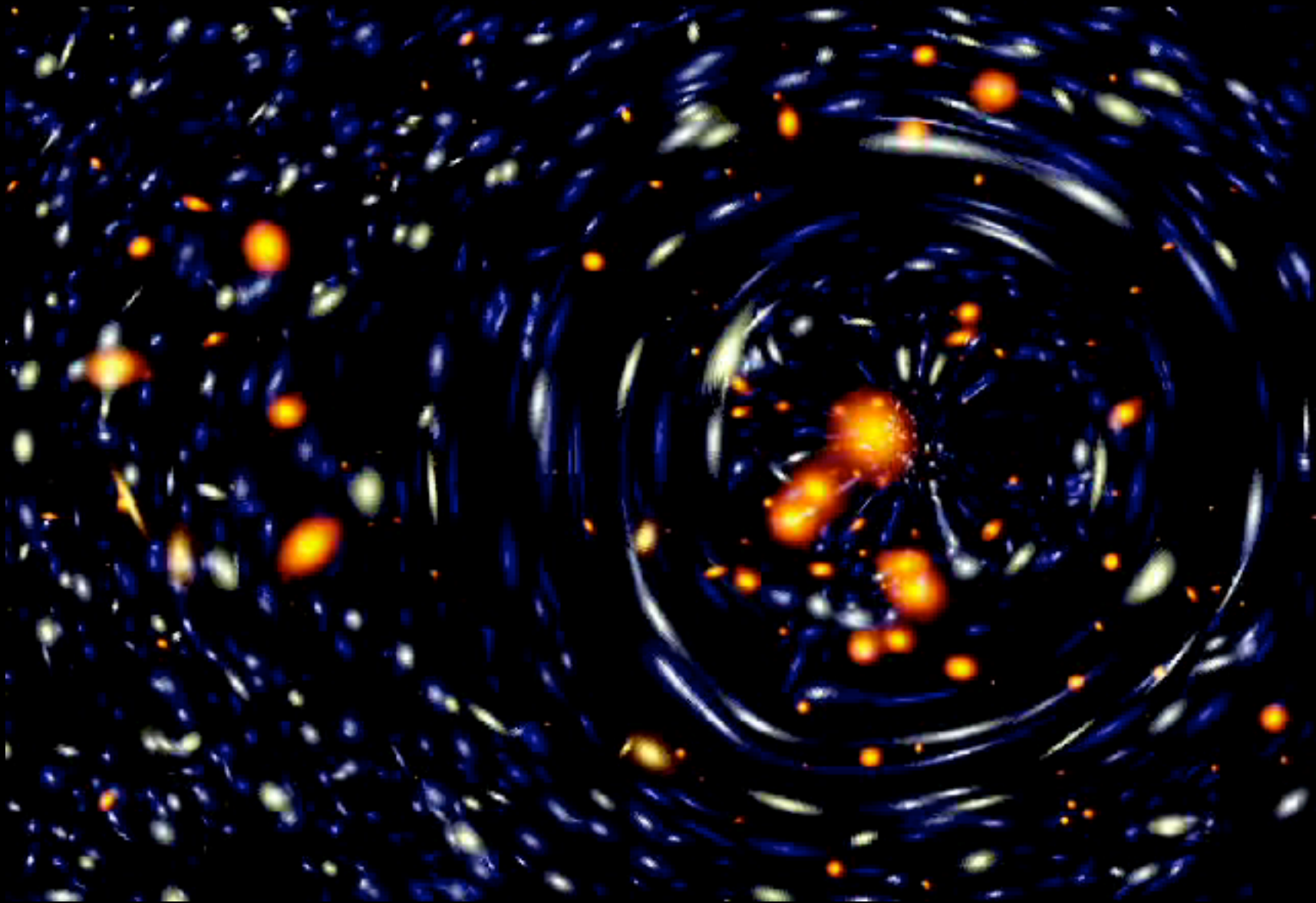
Searching for strong gravitational lenses





Einstein Ring Gravitational Lenses
Hubble Space Telescope • Advanced Camera for Surveys

Gravitational lens effect by a massive cluster



The size of the Einstein radius can be measured directly from the positions of the background galaxies.

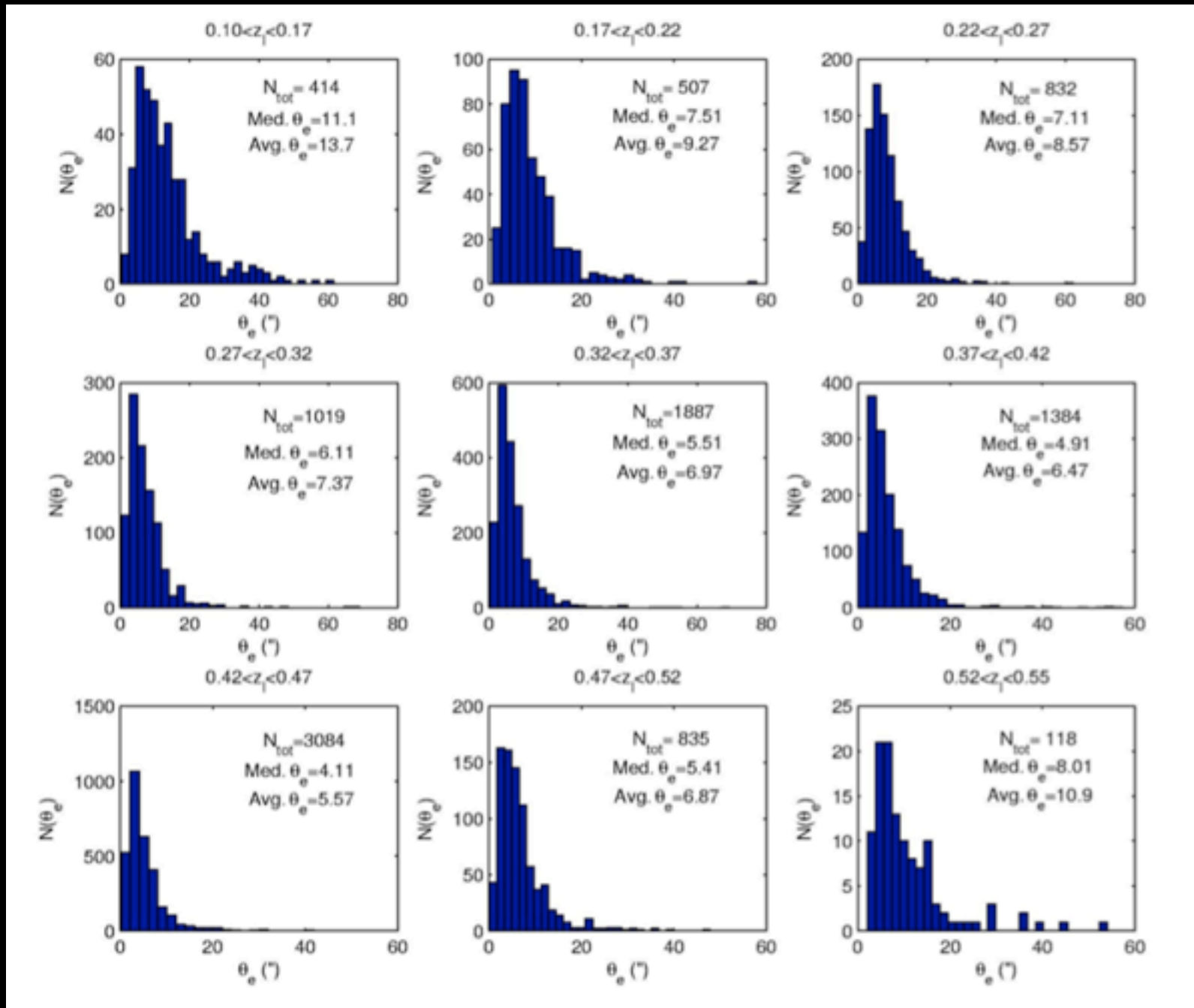
It can also provide us with the depth of the potential well (i.e. dark matter content) as well as the cosmological parameters of a given metric if the redshifts of the background galaxies can also be measured:

$$\theta_E = 4\pi \left(\frac{\sigma_{SIS}}{c} \right)^2 \frac{D_{LS}}{D_{OS}}$$

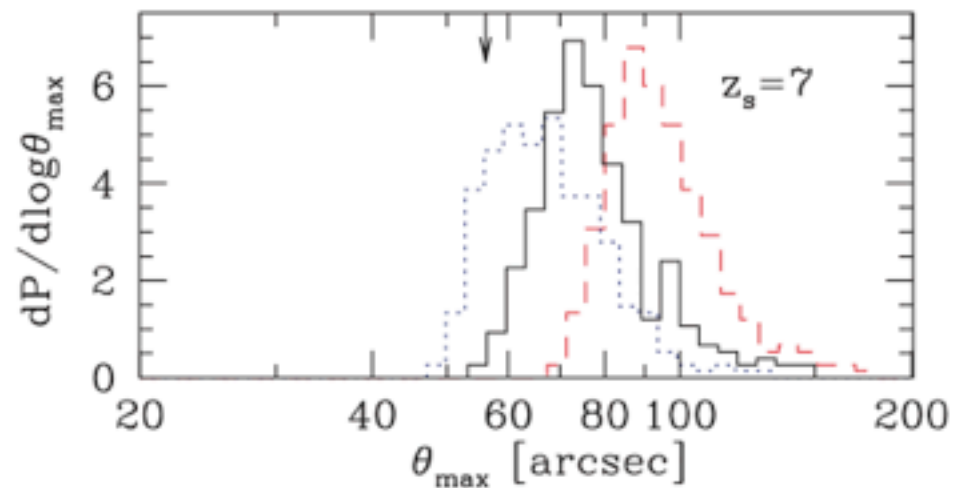
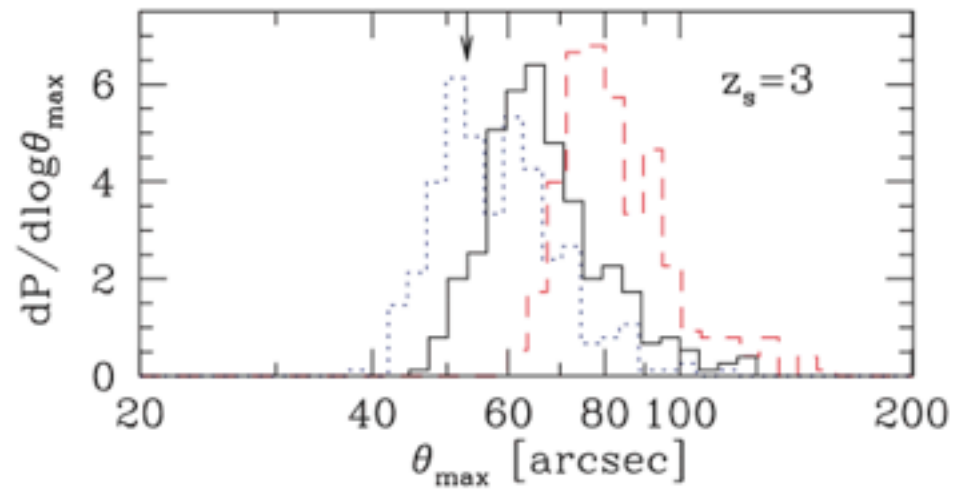
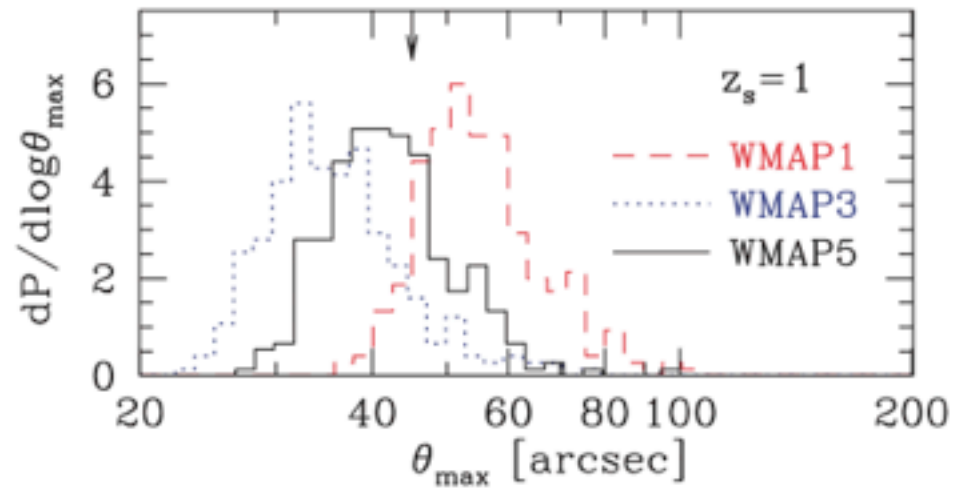
The diagram illustrates the relationship between the Einstein radius (θ_E), the spectrum of the lens (σ_{SIS}), and the spectrum of the source (D_{LS}/D_{OS}). The equation $\theta_E = 4\pi \left(\frac{\sigma_{SIS}}{c} \right)^2 \frac{D_{LS}}{D_{OS}}$ is shown. Three arrows point from the text labels below to the corresponding terms in the equation: one from 'positions of the images' to θ_E , one from 'spectrum of the lens' to $\left(\frac{\sigma_{SIS}}{c} \right)^2$, and one from 'spectrum of the source' to $\frac{D_{LS}}{D_{OS}}$.

positions of the images spectrum of the lens spectrum of the source

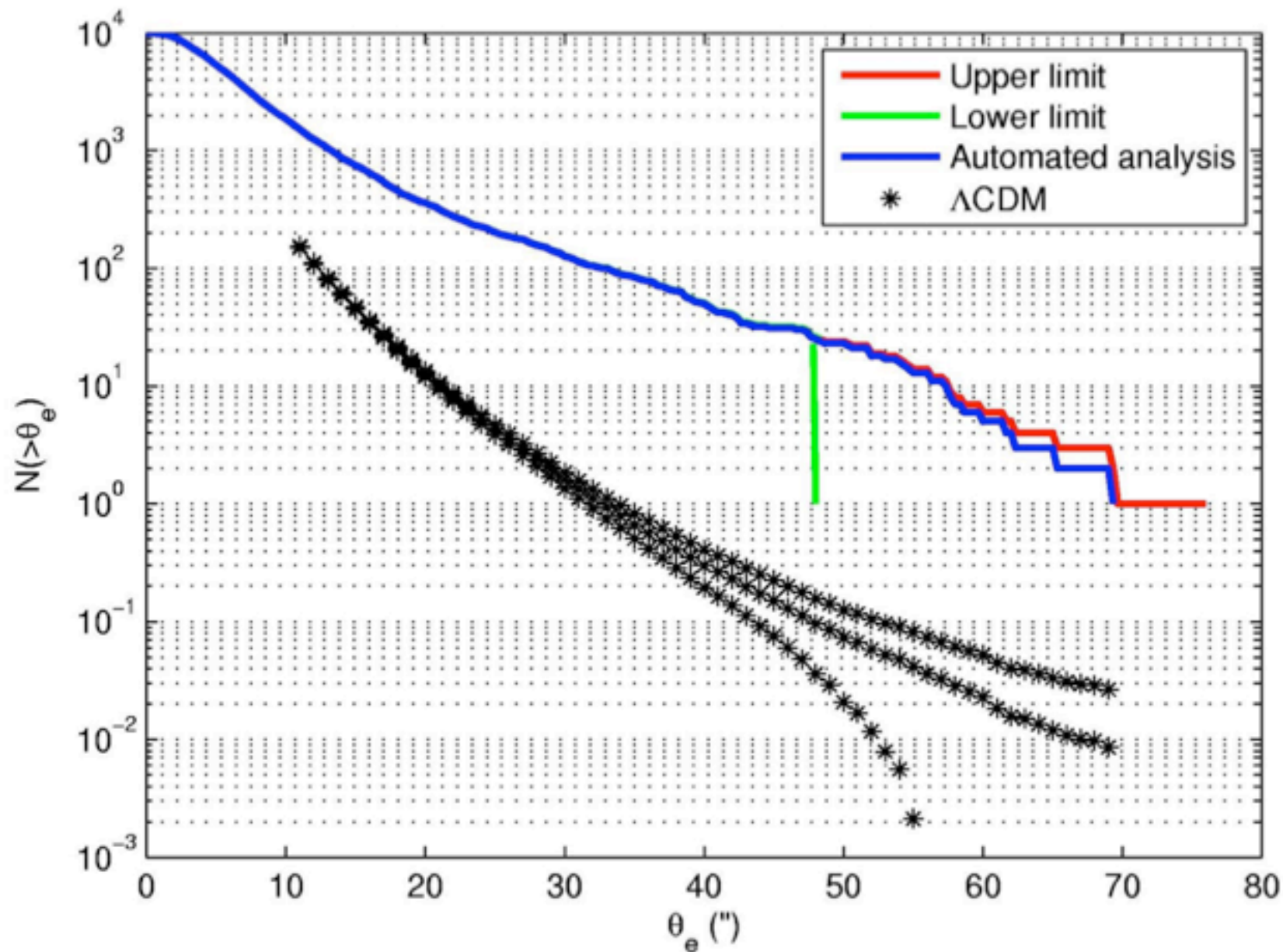
Inferring Einstein radii for 10,000 SDSS clusters

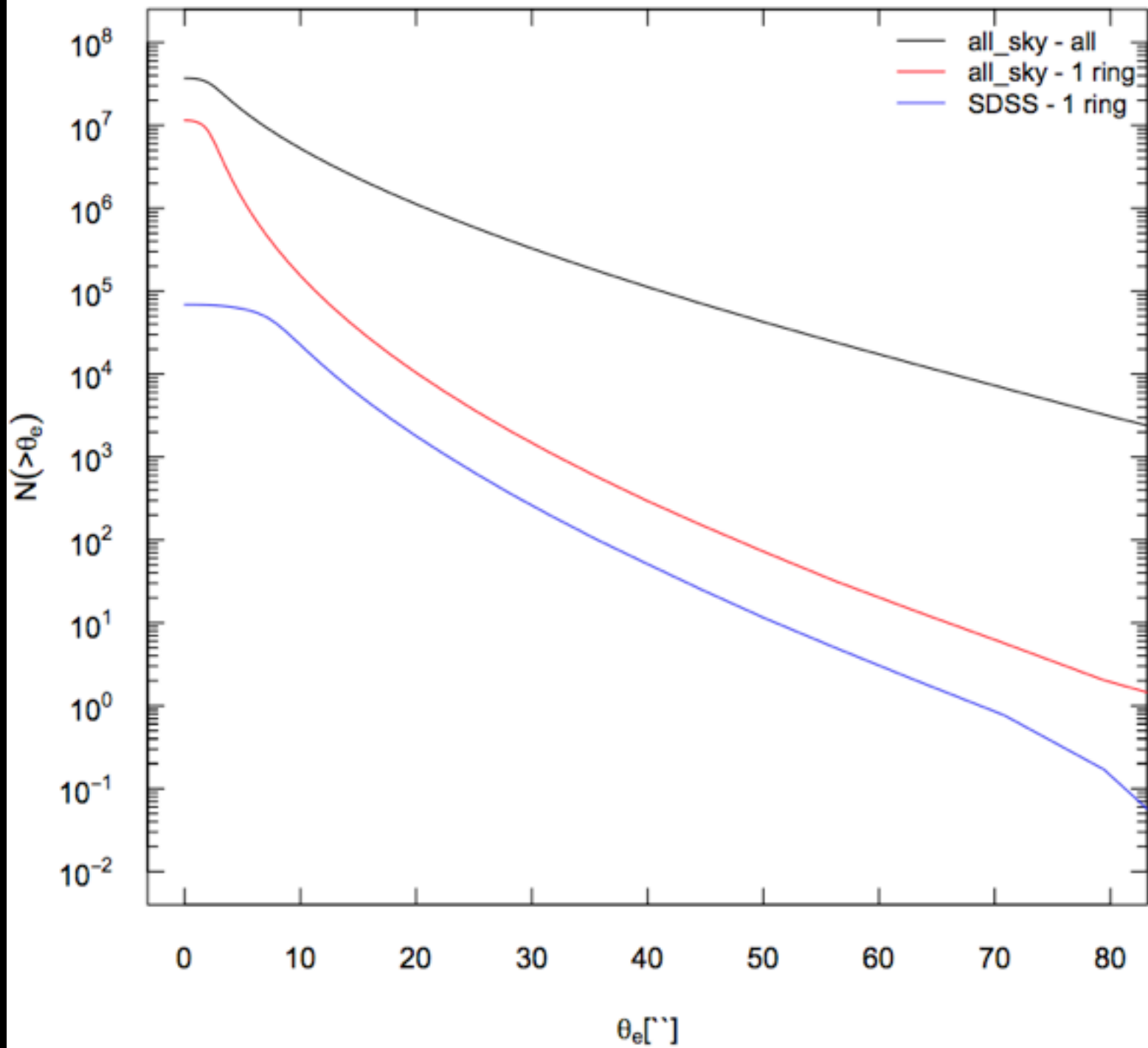


Theoretical expectations



Tension with Λ CDM scenario?





Summary

Methodology

- Proper statistical theory for extreme events
- Bayesian thresholding estimates work best
- Wide applications to astrophysics and cosmology

Limitations

- Assessment of selection biases in samples
- May require extensive testing with simulations