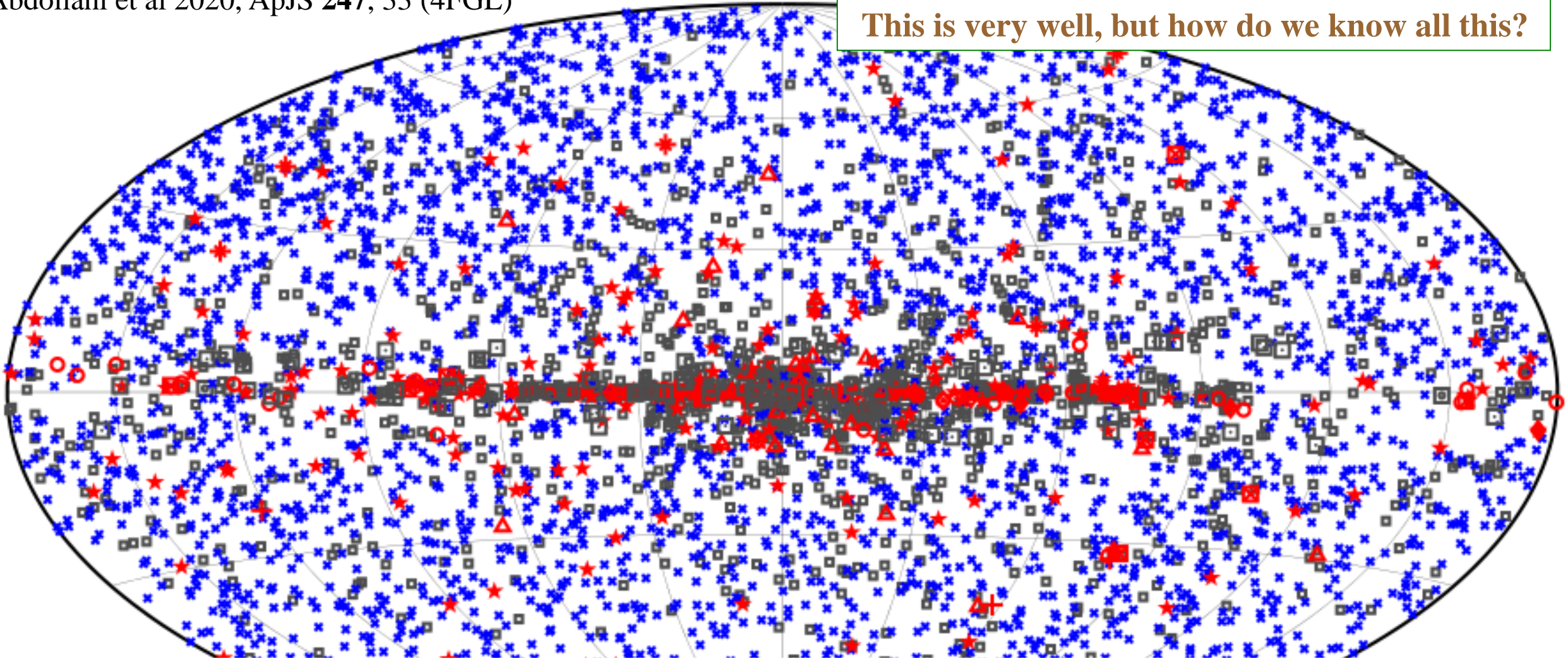


# Source association

1. Spatial association
2. Purity vs completeness
3. Chasing systematics
4. Other criteria
5. Galactic complications
6. Extended sources

This is very well, but how do we know all this?



□ No association	▣ Possible association with SNR or PWN	× AGN
★ Pulsar	△ Globular cluster	✳ Starburst Galaxy
⊠ Binary	+ Galaxy	○ SNR
★ Star-forming region	□ Unclassified source	◆ PWN
		★ Nova

# What problem do we want to address?

When a  $\gamma$ -ray source is found by chance, **how do we associate it with what we know from other wavelengths?**

- Applies to surveys (Fermi-LAT, eROSITA, CTA Galactic and extragalactic surveys)
- Critical for population studies and physical modelling
- Probabilistic approach
- Long history, first concepts date back to the 1970s for first radio surveys
- Often called **cross-match** in the literature
- Some concepts can apply to pointed observations when you want to assess the probability that what you see comes from a more common source class (*e.g.* blazars) than what you are looking for
- Of little use for extended sources, unfortunately

# What information do we need?

## What quantities do we expect will matter to this problem?

1. How well we localized the  $\gamma$ -ray source (the localization precision). In other contexts the localization precision of the counterparts may matter too (assume negligible here)
2. How many potential counterparts we consider (the counterpart density)
3. The plausibility that those counterparts emit  $\gamma$ -rays (not the same for stars and blazars). If possible, this is handled before, by selecting classes of sources that we **know** collectively emit  $\gamma$ -rays (blazars, pulsars)
4. The individual properties of the counterparts (flux, spectrum, ...)

Let us put that into equations

# Probabilistic framework

**We want to compare two hypotheses:**

1.  $H_0$ : A putative counterpart is close to a  $\gamma$ -ray peak by chance
2.  $H_1$ : The putative counterpart is actually the same as the  $\gamma$ -ray source

We will adopt a Bayesian approach:  $\Pr\{\mathbf{M}|\mathbf{D}\} = \Pr\{\mathbf{M}\} \Pr\{\mathbf{D}|\mathbf{M}\} / \Pr\{\mathbf{D}\}$

where  $\Pr(\mathbf{D})$  is just a normalization constant

# How do we get the localization precision?

## The instrument's Point Spread Function (PSF) is the key ingredient

1. If the PSF is the same for all events (not energy dependent, in particular), with dispersion  $\sigma$  along one axis, then the dispersion of the average over  $N$  events is  $\sigma / \sqrt{N}$
2. For many counts, the compound localization will converge to a **Gaussian** (central-limit theorem) of the same dispersion
3. In general, (not the same PSF for all events) the localization precision will be obtained from the logLikelihood using Wilks' theorem. Assuming that the source is truly at position  $\mathbf{r}_T$ ,  $\Delta = 2 \ln(L_{\max}/L_T)$  is distributed as  $\chi^2(2 \text{ dof})$  when the 2D position is fitted to the data. Particularly simple  $F(\Delta) = 1 - \exp(-\Delta/2)$ .  
Related to  $TS = 2 \ln(L_{\max}/L_0)$  used for assessing the significance of a source

# How do we get the probability density under $H_1$ ?

## Definition of the localization error

1. Remember that  $F(\Delta) = 1 - \exp(-\Delta/2)$  with  $\Delta = 2 \ln(L_{\max}/L_T)$
2. The likelihood contours are not necessarily symmetrical (either due to instrumental characteristics or to background features such as other nearby sources) but again for enough counts the tip converges to a Gaussian propto  $\exp(-(\mathbf{r}/\sigma)^2/2)$ ,  $\ln L$  becomes a 2D paraboloid and the contours converge to **ellipses**.  
95% confidence contours:  $\Delta = -2 \ln(0.05) = -(\mathbf{R}_{95}/\sigma)^2$  so  $\mathbf{R}_{95} = \sqrt{-2 \ln(0.05)} \sigma \approx 2.45 \sigma$
3. Under  $H_1$ , in the simplest case of an error circle, the probability density of the distance between the  $\gamma$ -ray peak and the counterpart is  $f_T(\mathbf{r}) = r/\sigma^2 \exp(-(\mathbf{r}/\sigma)^2/2)$   
 $\mathbf{r} = \|\mathbf{r}_P - \mathbf{r}_T\|$  is viewed as a random variable in  $\mathbf{r}_P$  ( $\gamma$ -ray peak when the counterpart is known)

We neglect here complications related to the sphericity of the sky



## How do we get the probability density under $H_0$ ?

### In general, we start from a catalog of counterparts

1. Under  $H_0$ , if the counterpart density  $\rho$  is reasonably constant (for example AGN outside the Galactic plane), then, noting  $r$  the 2D distance to any point in the sky,  $dN/dr = 2\pi r\rho$
2. As long as all counterparts are considered equal, we will consider the **nearest one**
3. The probability of finding the nearest neighbor at  $\mathbf{x}$  is  $p(\mathbf{x}) = \Pr\{N(\mathbf{r}<\mathbf{x})=0\}$ . We can write  $p'(\mathbf{x}) = 2\pi x\rho p(\mathbf{x})$  so that  $p(\mathbf{x}) = \exp(-\pi x^2\rho)$
4. Under  $H_0$ , the probability density of the distance to the closest counterpart is  $f_R(r) = -p'(r)$  so that  $f_R(r) = 2\pi r\rho \exp(-\pi r^2\rho)$
5. To get there, we need the catalog to be complete (at a given flux limit). If the detection rate varies over the sky (*e.g.* AGN through the Galactic plane), it must be accounted for.

We assume in the following that the local counterpart density  $\rho$  can be obtained



# Likelihood ratio

**We compare the two probability densities (random and true)**

1.  $H_0: f_R(\mathbf{r}) = 2\pi r \rho \exp(-\pi r^2 \rho)$

2.  $H_1: f_T(\mathbf{r}) = r/\sigma^2 \exp(-(r/\sigma)^2/2)$

3. The likelihood ratio is  $LR(r) = \frac{f_T(r)}{f_R(r)} = \frac{1}{2\pi \rho \sigma^2} \exp\left(\pi \rho r^2 - \frac{r^2}{2\sigma^2}\right)$

4. There is no free parameter here:  $\sigma$  comes from the logLikelihood contours (specific to each source),  $\rho$  is assumed known (but can depend on direction in the sky) and  $r$  is simply the distance between the  $\gamma$ -ray peak and the counterpart (observed quantity).

5. No hope to ever find a reliable counterpart if  $K = 2\pi \rho \sigma^2 > 1$ . In that case (one random counterpart on average in  $R_{68}$ ) LR does not even decrease with  $r$ .

6. The likelihood ratio provides a ranking between associations (over a full catalog) but not a probability

## Association probability

**In the Bayesian approach, we must consider the a priori probabilities**

1. A priori  $\Pr\{\mathbf{H}_1\}$  and  $\Pr\{\mathbf{H}_0\}$  such that  $\Pr\{\mathbf{H}_1\} + \Pr\{\mathbf{H}_0\} = 1$  and  $\Gamma = \Pr\{\mathbf{H}_1\} / \Pr\{\mathbf{H}_0\}$
2. A posteriori  $\Pr\{H_1 | r\} = \frac{\Pr\{H_1\} f_T(r)}{\Pr\{H_1\} f_T(r) + \Pr\{H_0\} f_R(r)} = \frac{1}{1 + 1/(\Gamma LR(r))}$
3. At this point  $\Pr\{\mathbf{H}_1\}$  and  $\Pr\{\mathbf{H}_0\}$  (or  $\Gamma$ ) are not known yet.

Only spatial characteristics ( $\rho$  and  $\sigma$ ) are considered for now.

# Thresholding

## We want to define a probability threshold

1. A counterpart is considered safe if  $\Pr\{H_1|r\} = \frac{1}{1+1/(\Gamma \text{LR}(r))} > \beta$
2. Equivalent to  $\text{LR}(r) = \frac{1}{2\pi\rho\sigma^2} \exp\left(\pi\rho r^2 - \frac{r^2}{2\sigma^2}\right) > \frac{1}{\Gamma(1/\beta-1)} = \alpha$
3. Or to  $\frac{r}{\sigma} < \sqrt{-\frac{2 \ln(K\alpha)}{(1-K)}}$  where  $\mathbf{K} = 2\pi\rho\sigma^2$
4. No association can be found if  $\mathbf{K} > 1$  (already seen) or  $\mathbf{K} > 1/\alpha$ , if  $\rho$  is too large. Usually  $\alpha > 1$  so the second condition is more stringent.
5. The  $r/\sigma$  threshold depends only on  $\alpha$ , not separately on  $\Gamma$  and  $\beta$ . It decreases with  $\mathbf{K}$  and  $\alpha$  so sources will be accepted at larger  $r/\sigma$  for smaller  $\rho$ , larger  $\Gamma$  and smaller  $\beta$
6. Ex:  $\beta = 0.8$ ,  $\Gamma = 1/2$ ,  $\rho = 1/\text{sq deg} \rightarrow \alpha = 8$ . Association possible if  $\sigma < 0.141^\circ$ .

# Catalog of $\gamma$ -ray sources

## We are now considering entire catalogs

1. We work with a catalog of  $M$   $\gamma$ -ray sources indexed by  $i$ , with localization precision  $\sigma_i$
2. We note  $p_i = \Pr\{H_1 | r_i\} = \frac{1}{1+1/(\Gamma \text{LR}(r_i))}$ . Good association if  $p_i > \beta$
3. Remember that the threshold in  $r/\sigma$  depends only on  $\alpha$ , not separately on  $\Gamma$  and  $\beta$ . So we can decide that we will set  $\beta$  to 0.8, say (the same for all counterpart catalogs), and it remains to choose  $\Gamma$  (separately for each counterpart catalog).

The localization of counterparts is assumed to be better than the  $\gamma$ -ray localization (in general, a few arcsecs vs a few arcmins) and we consider only the **closest one** in this simple approach

# False associations

## How many false associations do we expect?

1. For each source such that  $p_i > \beta$ , the number of false associations is a random variable  $F_i$  whose value is either 0 or 1.
2. In this framework the distance  $r_i$  is not a random variable but an observed quantity and the localization error  $\sigma_i$  is known from the logL contours. So  $F_i$  is a simple Bernoulli variable with probability  $1 - p_i$ . Its expectation is  $E_T(F_i) = 1 - p_i$  and its variance is  $V_T(F_i) = p_i (1 - p_i)$
3. The total number of false associations is  $F = \sum\{p_i > \beta\} F_i$ . Its expectation is  $E_T(F) = \sum\{p_i > \beta\} (1 - p_i)$  and the sources are independent so its variance is  $V_T(F) = \sum\{p_i > \beta\} p_i (1 - p_i)$
4. By construction all  $p_i > \beta$  so  $\beta E_T(F) < V_T(F) < E_T(F)$ , close to the Poisson regime  $V(F) = E(F)$
5.  $p_i$  depends on the a priori probability ratio  $\Gamma$ , so the expected number of false associations is a function of  $\Gamma$ . When  $\Gamma \ll 1$  ( $H_1$  very unlikely),  $p_i < \beta$  for all sources so  $E_T(F) = 0$ . When  $\Gamma \gg 1$  ( $H_1$  very likely),  $p_i \rightarrow 1$  for all sources so  $E_T(F) = 0$  too.  $E_T(F)$  reaches a maximum for moderate  $\Gamma$

## False associations 2

### How many false associations do we expect?

1. We can also estimate the number of false associations in a different manner, either by simulations (move the  $\gamma$ -ray sources, apply the procedure and count) or by a simple surface estimate.
2. In this framework the distance  $\mathbf{r}_i$  is again the random variable in  $H_0$  so that the cumulative probability of having one random counterpart within  $\mathbf{r}_i$  is  $\mathbf{F}_R(\mathbf{r}_i) = 1 - \exp(-\pi \mathbf{r}_i^2 \rho)$ .

3. Accept associations up to  $r_i^{\max} = \sigma \sqrt{-\frac{2 \ln(K_i \alpha)}{(1-K_i)}}$  where  $\mathbf{K}_i = 2\pi \rho \sigma_i^2$  and  $\alpha = 1/(\Gamma(1/\beta-1))$

4. Expected number of false positives  $E_R(F_i) = F_R(r_i^{\max}) = 1 - \exp\left(\frac{K_i \ln(K_i \alpha)}{(1-K_i)}\right)$

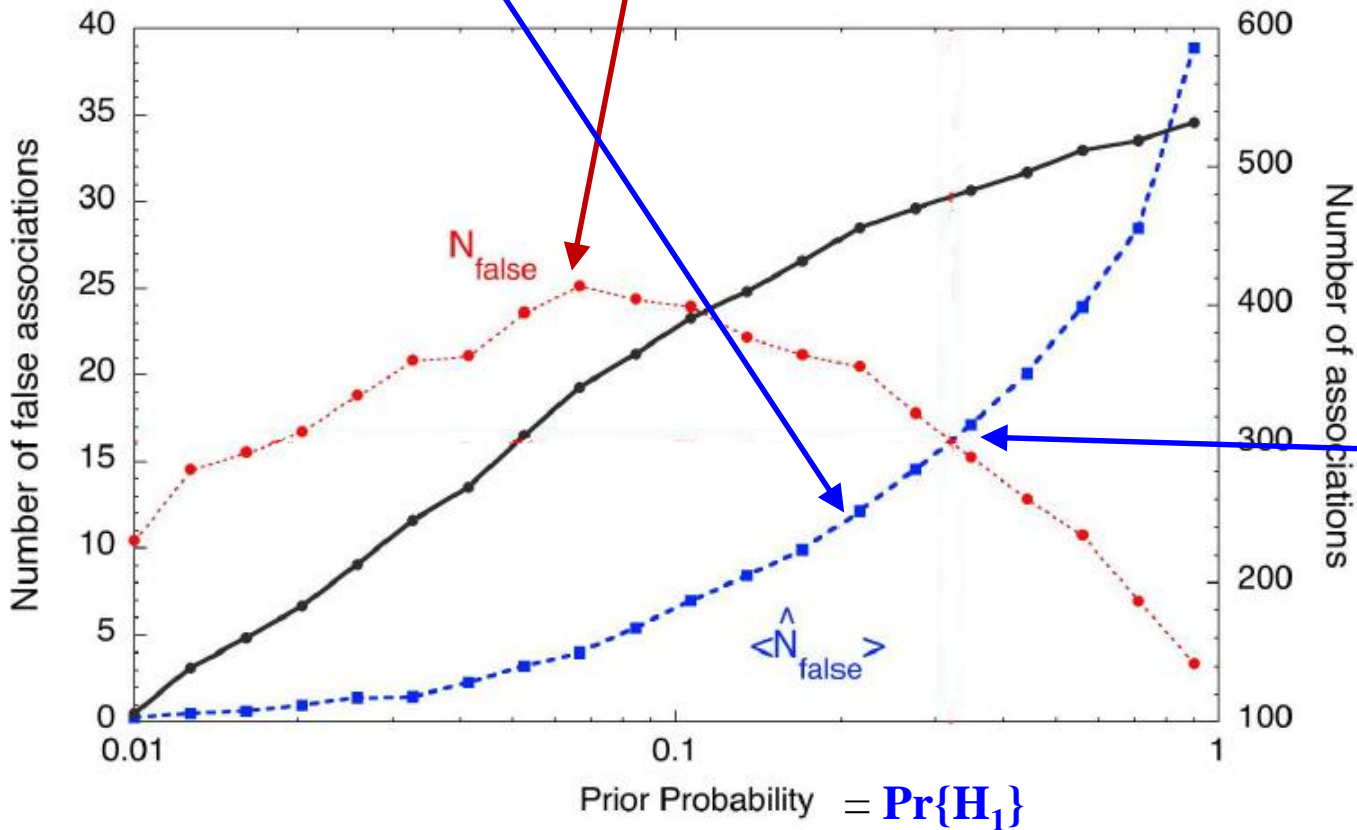
5. Summed over all sources  $\mathbf{E}_R(\mathbf{F}) = \Sigma \mathbf{F}_R(\mathbf{r}_i^{\max})$  restricted to  $\mathbf{K}_i < \min(1, 1/\alpha)$

6.  $\mathbf{K}_i$  does not depend on  $\Gamma$ , and  $\alpha$  decreases with  $\Gamma$ , so  $\mathbf{E}_R(\mathbf{F})$  increases with  $\Gamma$  from 0 to  $\mathbf{M}$

# Defining the prior probability

**We must reconcile the two estimates of false associations**

Writing  $E_R(F) = E_T(F)$  results in an equation over  $\Gamma$  or  $\Pr\{H_1\}$  that can be solved numerically



Example of such curves

Scale for false associations is at left

Black curve (with scale at right) is the total number of associations

The correct choice of  $\Pr\{H_1\}$  is where the red and blue curves intersect, at  $\Pr\{H_1\} \approx 1/3$  or  $\Gamma \approx 1/2$

This means that we expect about 1/3 of the  $\gamma$ -ray sources to be among those counterparts

The reliability  $R$  (called **precision** in statistics) is the fraction of true associations among accepted ones  $R = 1 - E_T(F)/N_{\text{assoc}}$



# True associations

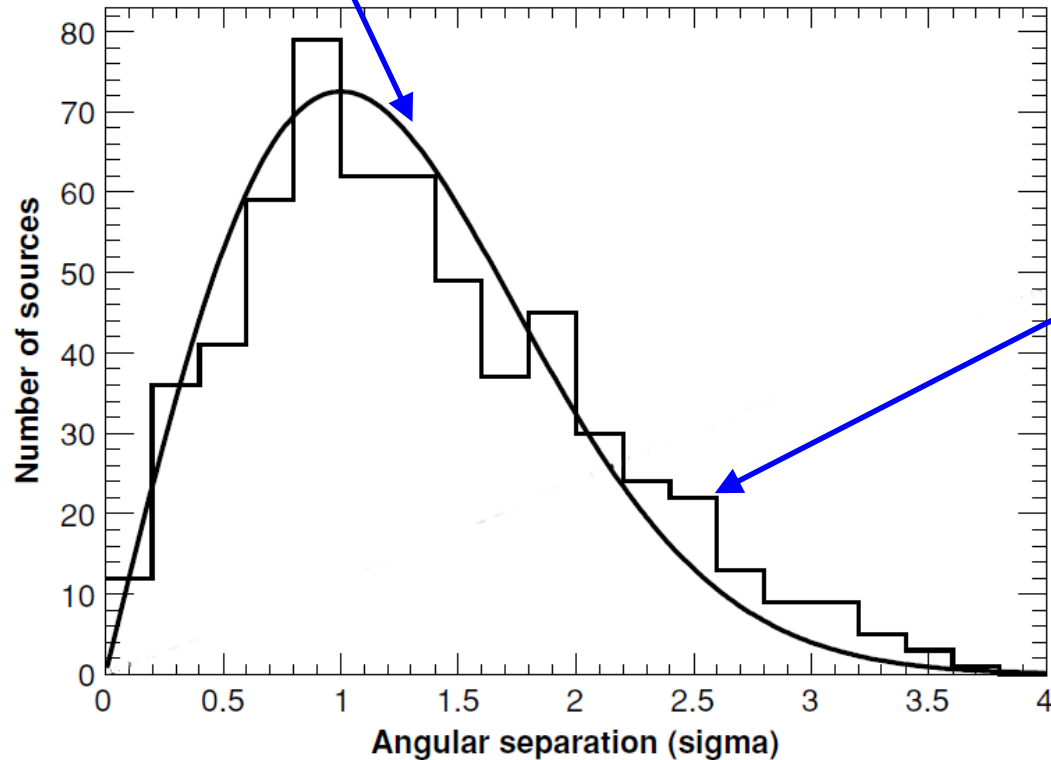
## How many true associations do we expect?

1. The number of true associations is also a random variable  $\mathbf{T}_i$  whose value is either 0 or 1.
2.  $\mathbf{T}_i$  is again a simple Bernoulli variable with probability  $\mathbf{p}_i$ . Its expectation is  $\mathbf{E}(\mathbf{T}_i) = \mathbf{p}_i$  and its variance is  $\mathbf{V}(\mathbf{T}_i) = \mathbf{p}_i (1 - \mathbf{p}_i)$
3. The total number of true associations among the accepted sources is  $\mathbf{T}_{\text{acc}} = \sum\{\mathbf{p}_i > \beta\} \mathbf{T}_i$ . Its expectation is  $\mathbf{E}(\mathbf{T}_{\text{acc}}) = \sum\{\mathbf{p}_i > \beta\} \mathbf{p}_i$  and its variance is  $\mathbf{V}(\mathbf{T}_{\text{acc}}) = \sum\{\mathbf{p}_i > \beta\} \mathbf{p}_i(1 - \mathbf{p}_i) = \mathbf{V}_T(\mathbf{F})$
4. By construction all  $\mathbf{p}_i > \beta$  so  $\mathbf{0} < \mathbf{V}(\mathbf{T}_{\text{acc}}) < (1 - \beta) \mathbf{E}(\mathbf{T}_{\text{acc}})$ , way below Poisson regime  $\mathbf{V}(\mathbf{F}) = \mathbf{E}(\mathbf{F})$
5. The total number of true associations is  $\mathbf{T}_{\text{tot}} = \sum \mathbf{T}_i$ . Its expectation is  $\mathbf{E}(\mathbf{T}_{\text{tot}}) = \sum \mathbf{p}_i$
6. The completeness  $\mathbf{C}$  (called **recall** in statistics) is the fraction of accepted true associations among all true associations and can be estimated as  $\mathbf{E}(\mathbf{T}_{\text{acc}}) / \mathbf{E}(\mathbf{T}_{\text{tot}})$

# The Rayleigh distribution

**We must quantify the quality of the procedure**

The distribution of distances  $\mathbf{r}$  differs for each source but that of  $\mathbf{r}/\sigma$  is always the same  
 $f_T(\mathbf{r}/\sigma) = x \exp(-x^2/2)$  : Rayleigh distribution (2D Gaussian in polar coordinates)



Example on real sources

The black histogram is the observed distribution of  $\mathbf{r}/\sigma$

The curve is the Rayleigh distribution normalized to the number of sources

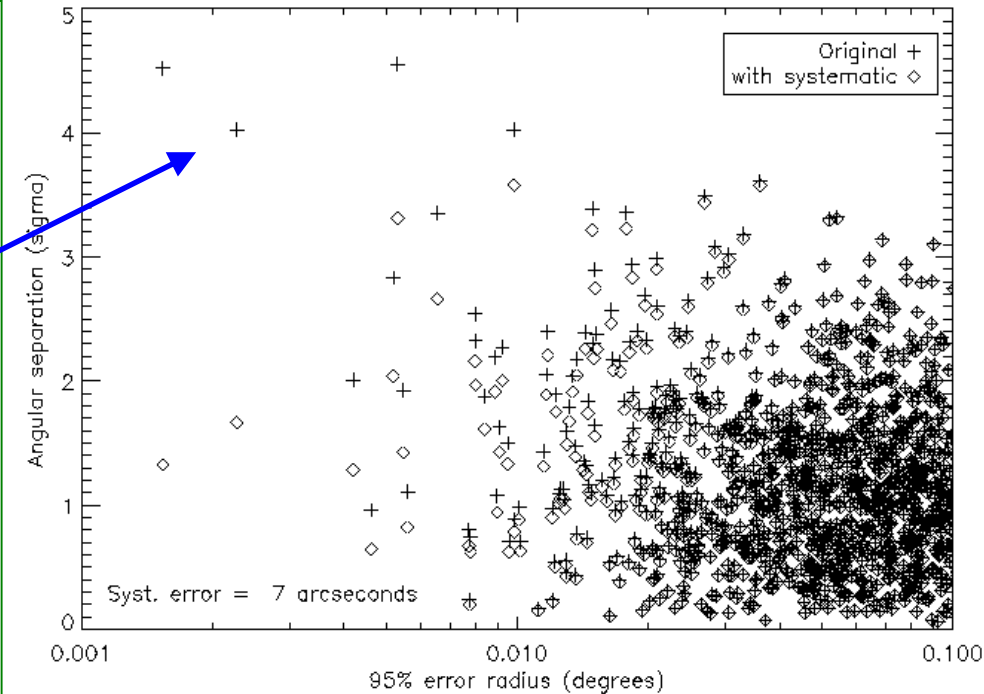
They can be compared by a Kolmogorov-Smirnov test

This example is not perfect. The histogram has a distinct tail, implying that something could be improved

# Chasing systematics

## What can explain a tail in the observed distribution of $r/\sigma$ ?

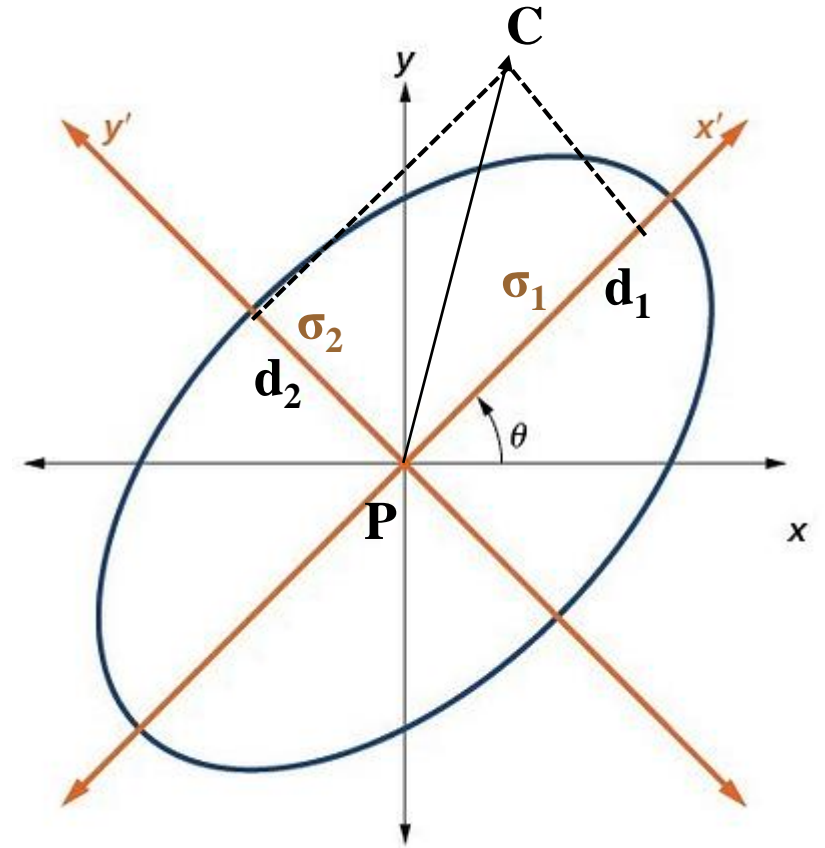
1. In general, the culprit is our estimate of the localization precision  $\sigma$
2. It can be wrong in two ways:
  - An absolute systematic error  $\sigma_{\text{abs}}$  (due to imperfect knowledge of the pointing direction) will affect the bright sources and can be checked by looking at bright known sources
  - A relative systematic error  $f_{\text{rel}}$  (due to confusion or background modeling) will affect all sources.  $f_{\text{rel}}$  can be fit to optimize the Rayleigh plot
3. Combined as  $\sigma_{\text{tot}}^2 = (f_{\text{rel}} \sigma)^2 + \sigma_{\text{abs}}^2$



# Elliptical errors

## How to go beyond a simple error circle?

1. In general, the localization region is an ellipse defined by two errors  $\sigma_1$  and  $\sigma_2$  and an angle  $\theta$  (can be wrt North or West)
2. In that case, the counterpart position  $\mathbf{C}$  wrt the  $\gamma$ -ray source  $\mathbf{P}$  must be expressed in the ellipse axes  $\rightarrow (\mathbf{d1}, \mathbf{d2})$
3. The previous formulae can then be used, replacing  $\sigma$  by  $\sqrt{(\sigma_1 \sigma_2)}$  and  $(\mathbf{r}/\sigma)^2$  by  $(\mathbf{d}_1/\sigma_1)^2 + (\mathbf{d}_2/\sigma_2)^2$



## Counterparts in large catalogs

### How can we do when the counterpart density is too large?

1. An often-used method is to consider the counterpart flux  $S_k$
2. The idea then is to define the source density only from those sources with flux no less than  $S_k$   
 $\rho_k = N(S \geq S_k) / \Omega$  (it should ideally be differential at  $S_k$ , not addressed here)
3. The likelihood ratio can then be expressed, replacing  $\rho$  by  $\rho_k$ , except that we will now check all possible pairs,  $LR_{ik} = \exp(-(r_{ik}/\sigma_i)^2/2) / (2\pi\rho_k\sigma_i^2)$  (no nearest neighbor term in exp)
4. We will then consider for each  $\gamma$ -ray source the counterpart with the largest likelihood ratio instead of the nearest neighbor:  $LR_i = \max_k LR_{ik}$
5. It is not so easy to define a global a priori probability  $\Pr\{H_0\}$

# Reliability in likelihood ratio method

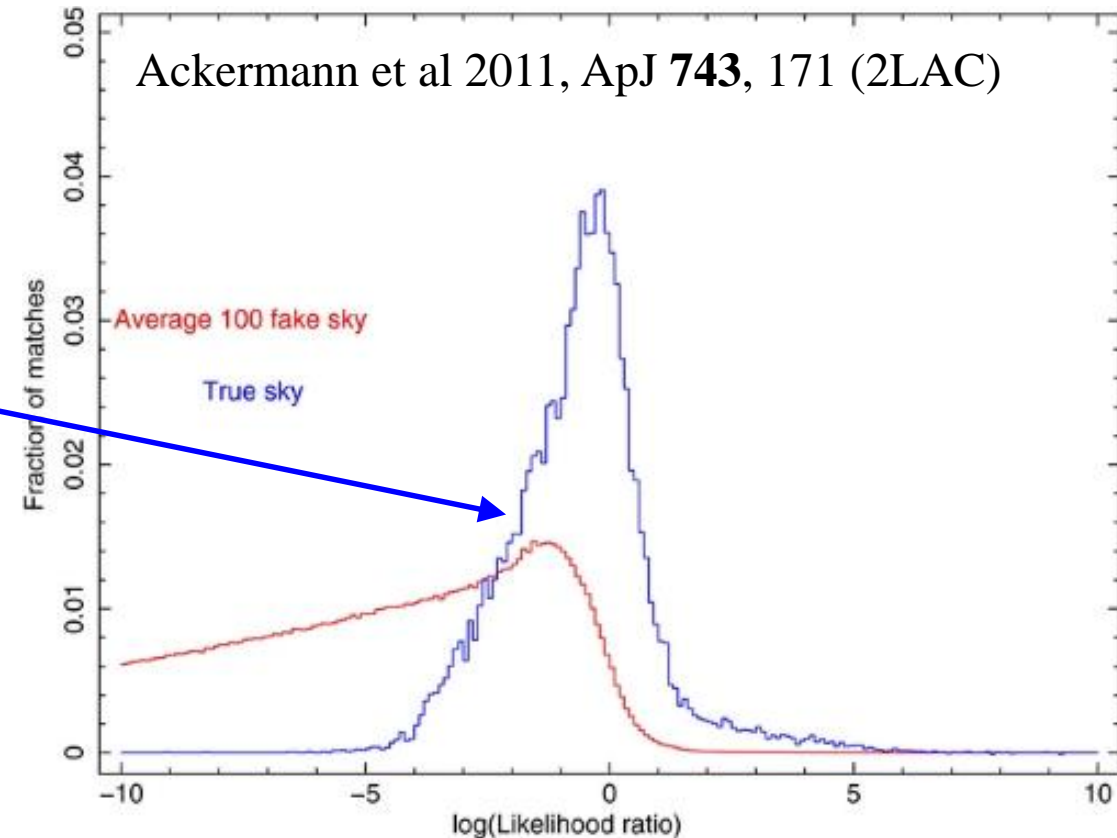
## How can we go further?

1. We can try to estimate the distribution of LR under  $H_0$  by simulating  $\gamma$ -ray sources randomly (but with a similar spatial distribution)

2. It is then possible to obtain from the true and the random **LR** distributions a reliability. Noting  $N_T$  and  $N_R$  the numbers of sources in a given **log(LR)** bin we define

$$R(LR) = \frac{1}{1 + N_R(LR)/N_T(LR)}$$

3. This is an approximate probability of association. It is noisy (because  $N_T$  is noisy) so it must be approximated by some analytic function.



## Adding other criteria

### Can we use other counterpart characteristics?

1. Yes. In the likelihood ratio method we can multiply the spatial term **S** by other terms (other data)
2. In the Bayesian formalism,  $\Pr\{\mathbf{D}|\mathbf{M}\} = \Pr\{\mathbf{S}|\mathbf{M}\} \times \Pr\{\mathbf{C}|\mathbf{M}\}$  (actually probability densities)
3. We must know the distributions of the secondary quantity **C** under  $\mathbf{H}_0$  and  $\mathbf{H}_1$
4. The distribution of **C** under  $\mathbf{H}_0$  is taken from that of the full counterpart catalog
5. The distribution of **C** under  $\mathbf{H}_1$  is taken from a subset of “sure” identifications (not so easy)
6. Can be easily generalized to multiple characteristics



# Complications

## Modern association tools

1. The localization precision of the counterparts must be accounted for (symmetric formulation)
2. Several counterpart catalogs must be handled together
3. The counterparts must be associated between themselves (in general they are better localized so we know whether an X-ray source is the same as an optical source or not)
4. Some sources can be absent at a particular wavelength simply because this source class emits little there (eg pulsars in the optical)

Implemented in the [NWAY package](#) by Mara Salvato et al (developed for eROSITA)

# Handling Galactic sources

## The Galactic plane is much more complex

1. The number of potential classes of emitters is much larger
2. The density of counterparts varies (latitude, distance to GC, spiral arms) so it must be estimated locally, for example via [kernel density estimation](#)
3. When the density of counterparts becomes too large locally (as explained before), it can become advantageous to consider them **collectively** (for example, star-forming regions rather than individual young stars, globular clusters rather than individual MSPs)
4. Galactic absorption/extinction biases the counterpart catalogs at many wavelengths (soft X-rays, optical/UV) whereas  $\gamma$ -ray sources are negligibly absorbed
5. Because of all that, care must be taken when simulating fake  $\gamma$ -ray sources to preserve their spatial distribution

# Handling extended sources

## Extended sources cannot fit into this probabilistic framework

1. Extended here means the radius **R is larger than the localization error** (not only larger than the PSF). Can be the same source even though the centroid is a little off, because the relative weights of emission regions inside the extended source can differ at different wavelengths
2. Many more parameters come into play: counterpart size (PWN, SNR)  $R_{\text{ctpt}}$ , the  $\gamma$ -ray size (when it can be measured)  $R_{\text{gam}}$ , the  $\gamma$ -ray localization (as before)  $R_{95}$  (95%)
3. **Majority** of Galactic TeV sources, unfortunately
4. When you find an extended source, look at images first! You can get images of the sky at many wavelengths from [NASA's SkyView](#).
5. The [Manitoba catalog](#) (Ferrand & Safi-Harb 2012, AdSpR **49**, 1313) is a good resource to look into

# Handling extended sources 2

Can we attempt to quantify this anyway?

1. Yes, sort of

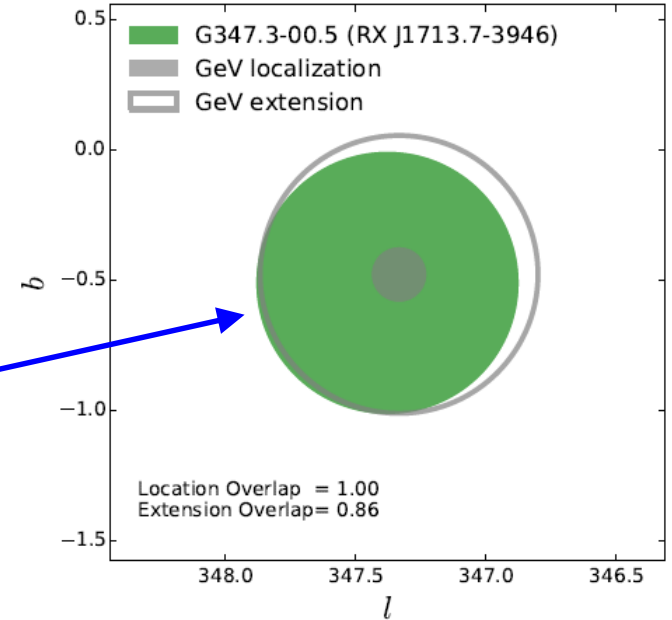
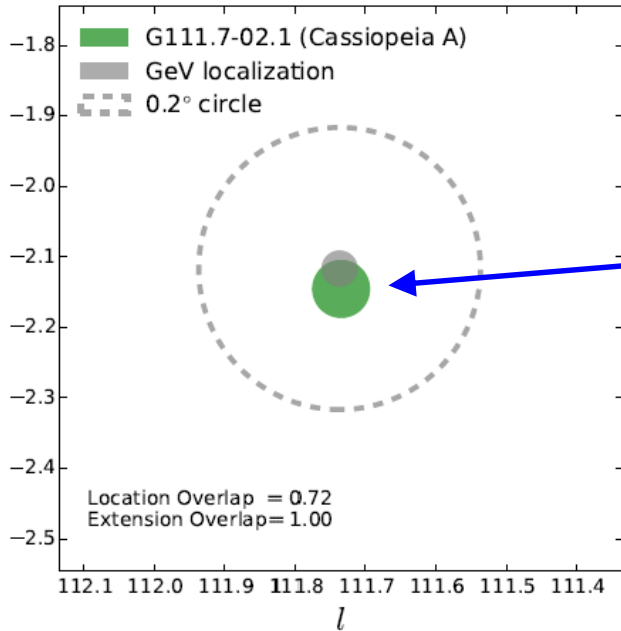
2. Location:  $O_{loc} = \frac{S_{ctpt} \cap S_{95}}{\min(S_{ctpt}, S_{95})} > O_{min}$

3. Extension:  $O_{ext} = \frac{S_{ctpt} \cap S_{gam}}{\min(S_{ctpt}, S_{gam})} > O_{min}$

4.  $O_{min}$  set to 0.5 or so

5. Estimate reliability from simulations

Acero et al 2016, ApJS **224**, 8 (SNRCat)



## References and credits

### Credits:

Jürgen Knödlseher (IRAP Toulouse) who created the Bayesian associations machinery for Fermi-LAT  
Benoit Lott (LP2i Bordeaux) who took up the task in 2015  
Dario Gasparrini (SSDC Roma) who handles the Likelihood Ratio machinery for Fermi-LAT

### References:

De Ruiter & Willis 1977 (A&AS **28**, 211) in the radio  
Sutherland & Saunders 1992 (MNRAS **259**, 413) for Likelihood Ratio  
Mattox et al 1997 (ApJ **481**, 95) for EGRET  
Rutledge et al 2000 (ApJS **131**, 335) for ROSAT  
Budavari & Szalay 2008 (ApJ **679**, 301) for combining probabilities over multiple catalogs  
Pineau et al 2017 (A&A 597, A89) for missing entries in some catalogs  
Salvato et al 2018 (MNRAS 473, 4937) introducing NWAY for eROSITA

# Conclusions

1. Well established framework for point sources
2. Statistical estimate of false and missed associations
3. Can accommodate source density, flux, other quantities
4. Extended sources more uncertain