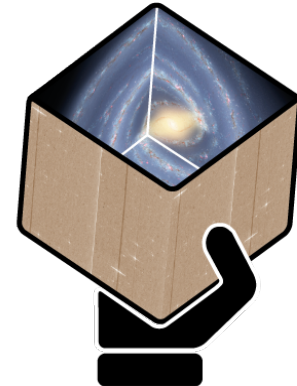


Working with astrometric data - warnings and caveats -

U. Bastian / X. Luri



gaia



Scientist's dream

- **Error-free data**
 - No random errors
 - No biases
 - No correlations
- **Complete sample**
 - No censorships
- **Direct measurements**
 - No transformations
 - No assumptions

Never ever available

Errors 1: biases

Bias:

your measurement is systematically too large or too small

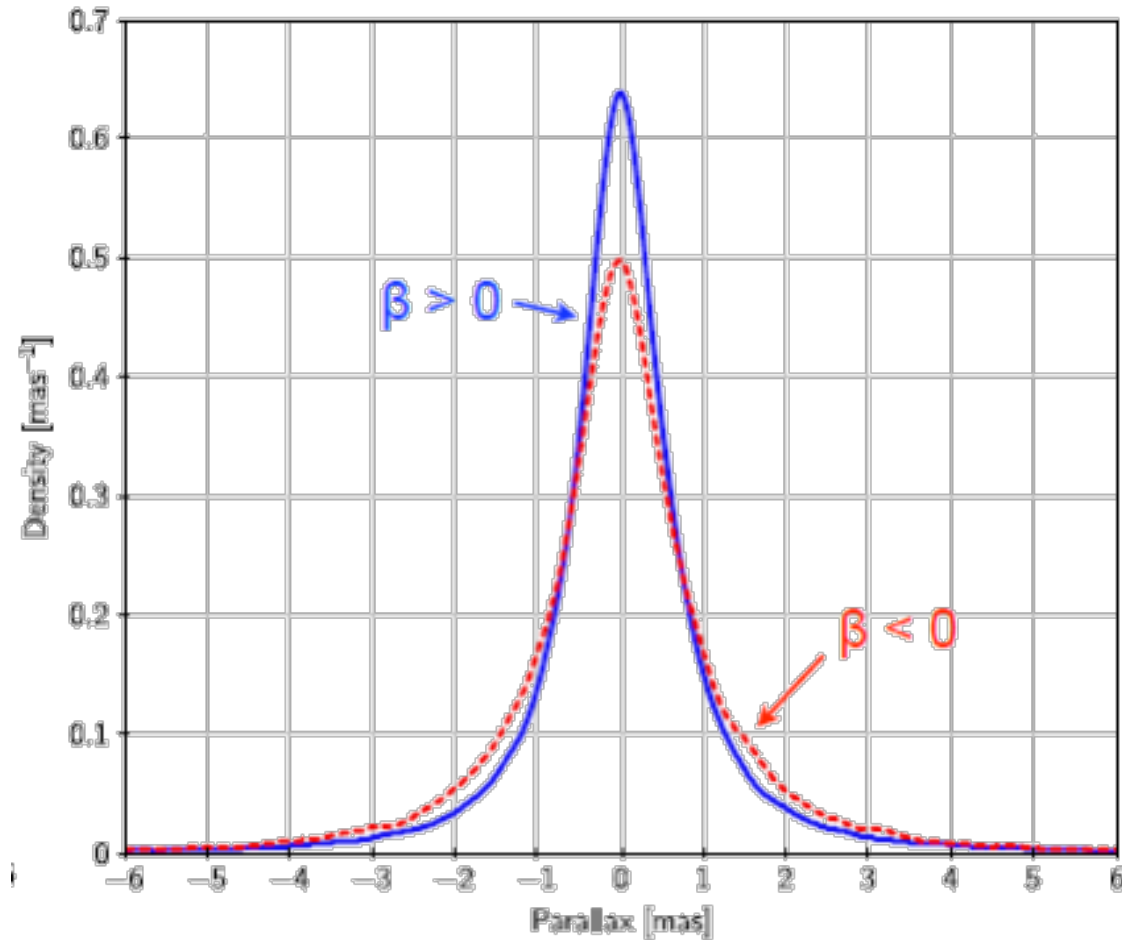
Example: DR1 parallaxes

- Probable global zero-point offset present; -0.04 mas found during validation
- Colour dependent and spatially correlated systematic errors at the level of 0.2 mas
- Over large spatial scales, the parallax zero-point variations reach an amplitude of 0.3 mas
- Over a few smaller areas (2 degree radius), larger parallax biases may occur of up to 1 mas

This is possibly the sole aspect in which Gaia DR1 is not better than Hipparcos (apart from the incompleteness for the brightest stars)
But see the Pleiades discrepancy ...

Global zero point from QSO parallaxes

TC5 (with colour terms)



$\beta > 0$: med(ϖ) = +0.002 mas

$\beta < 0$: med(ϖ) = -0.020 mas



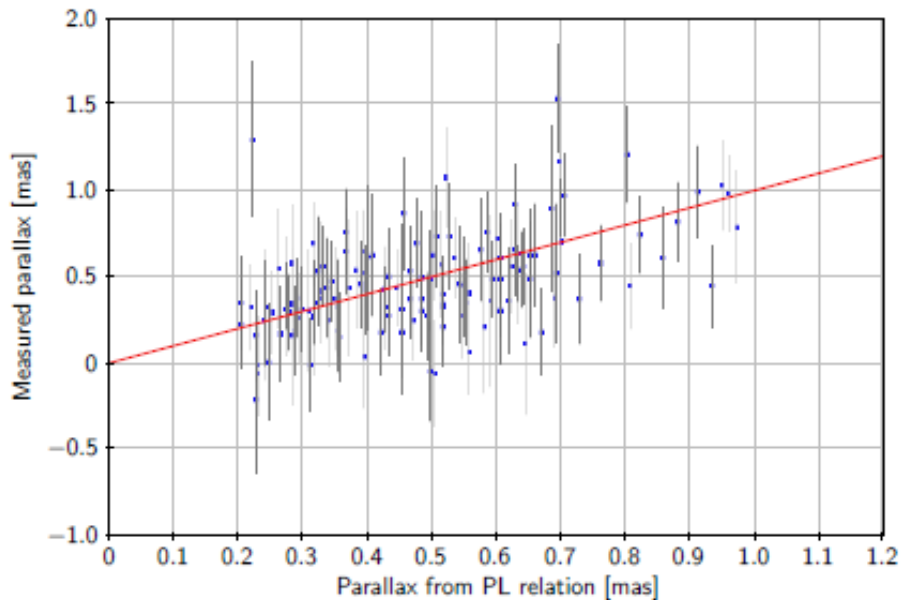
gaia



Global zero point from Cepheids

P-L relation from Tammann et al. (2003):

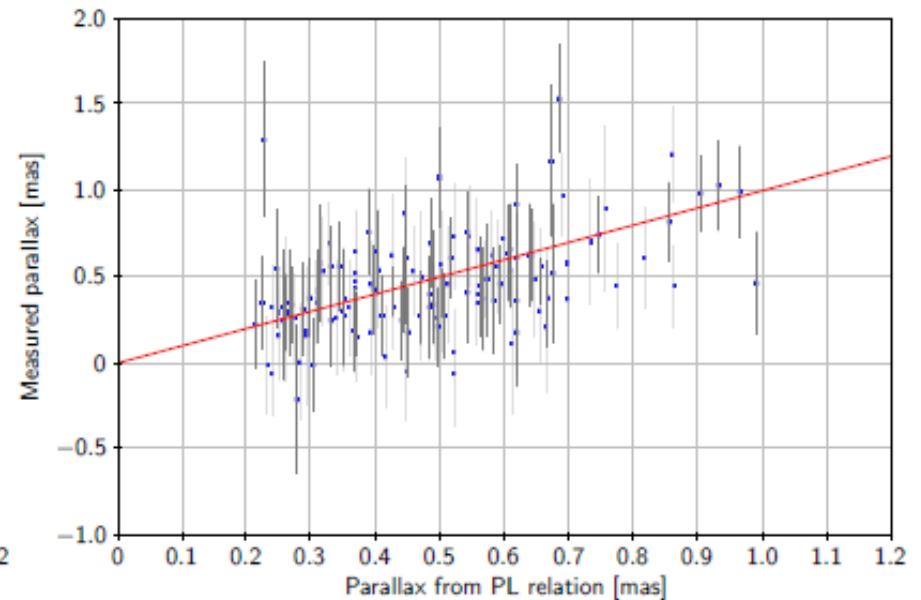
$$M_V = -3.141 \log P - 0.820$$



$$\text{med}(\Delta\varpi) = -0.015 \text{ mas}$$

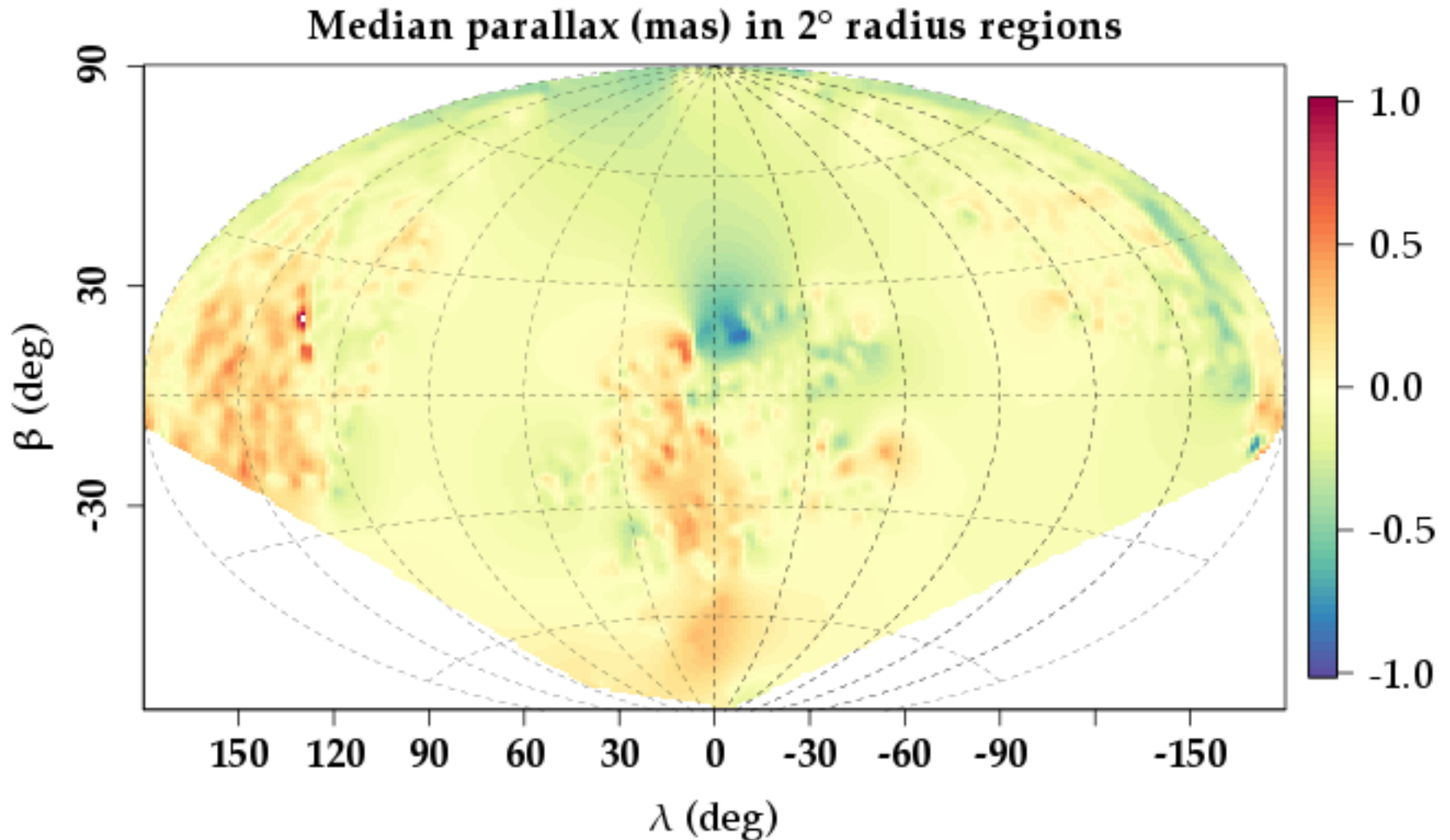
P-L relation from Fouqué et al. (2007)

$$M_V = -2.678 \log P - 1.275$$

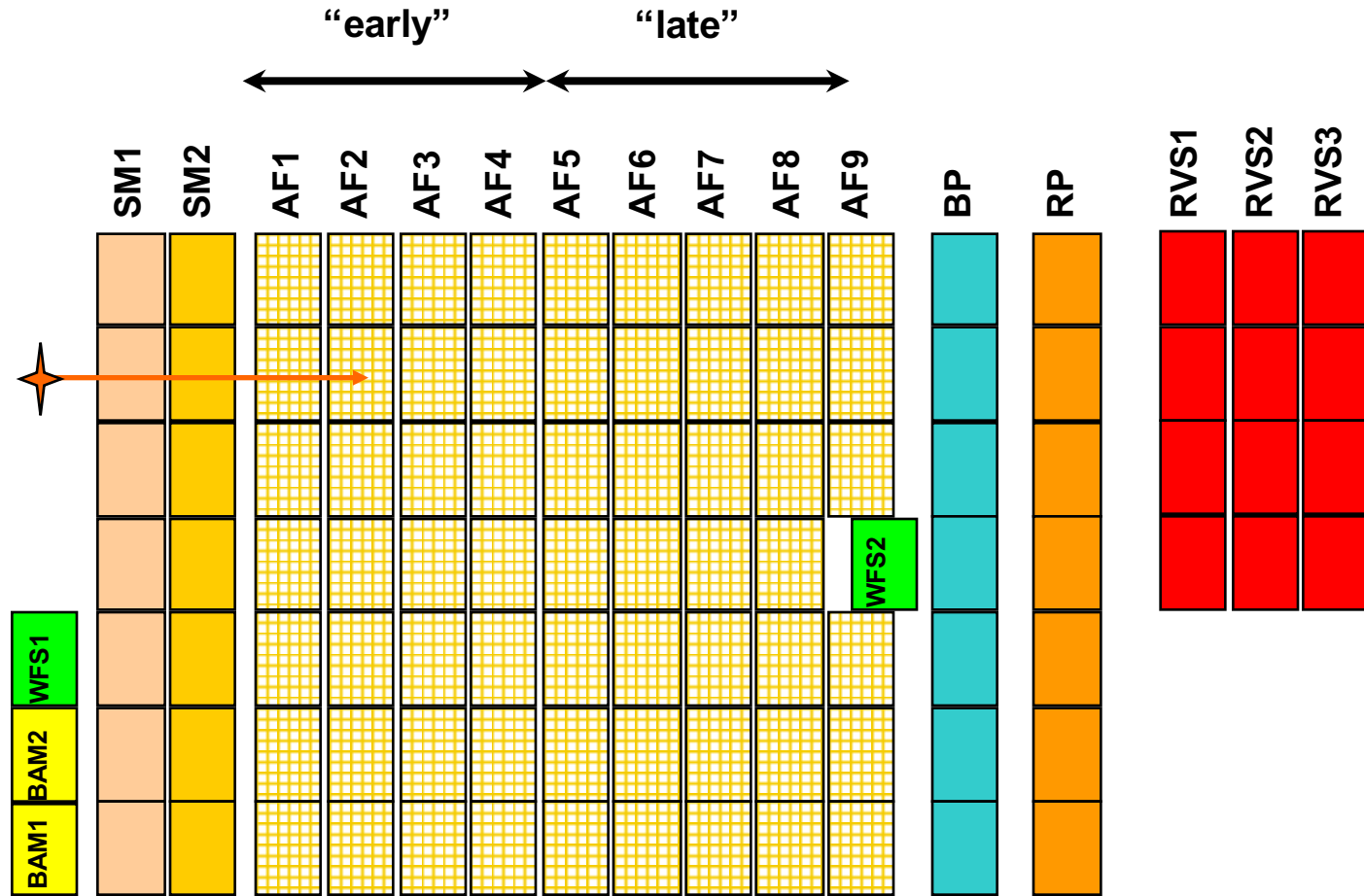


$$\text{med}(\Delta\varpi) = -0.017 \text{ mas}$$

Regional effects from QSOs (ecliptic coordinates)

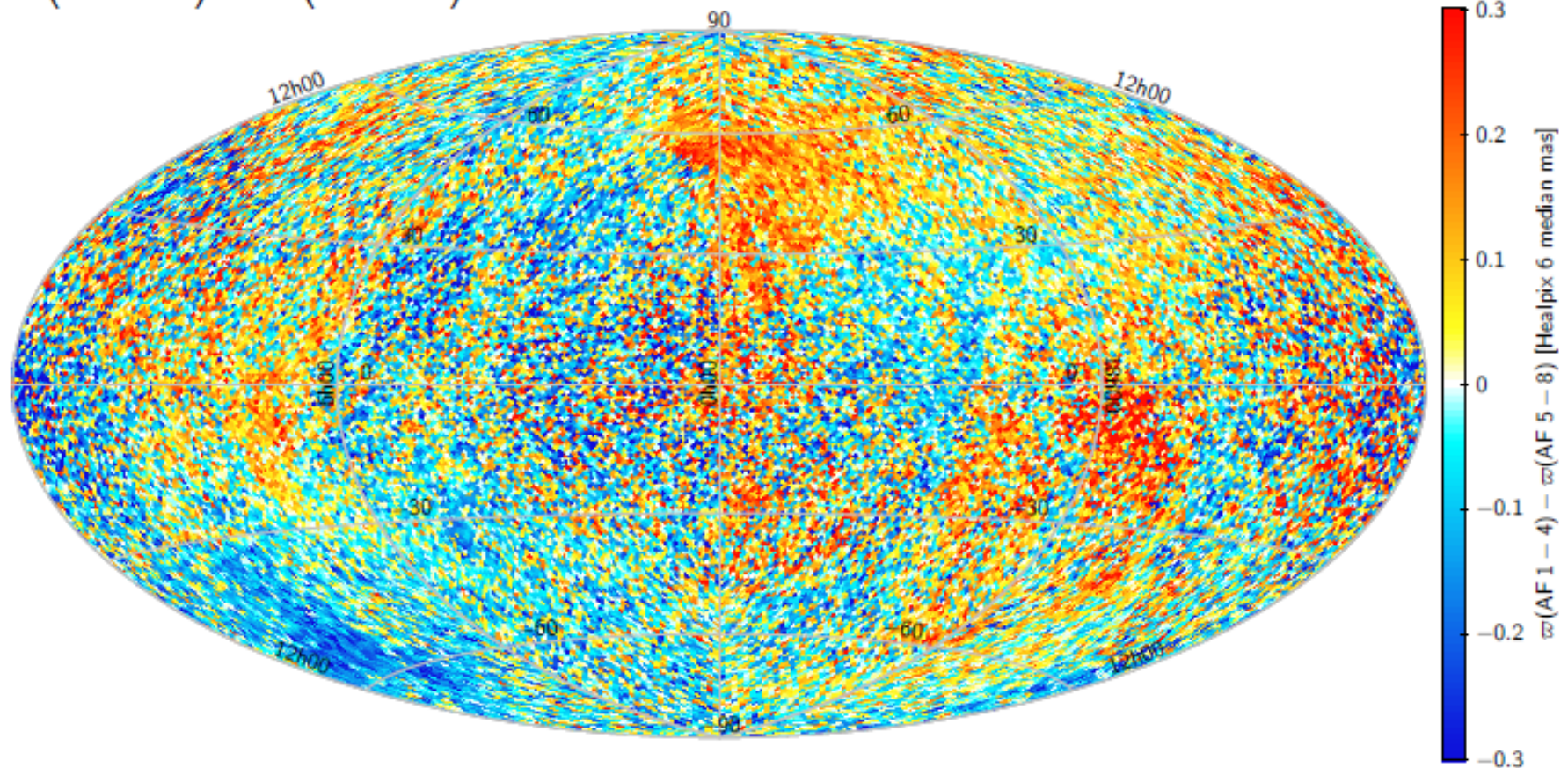


Split FoV



Regional effects from split FOV solutions (equatorial coordinates)

$$\varpi(\text{AF1-4}) - \varpi(\text{AF5-8})$$



Median value per pixel ($\sim 1 \text{ deg}^2$)

A. Bombrun

How to take this into account

- You can introduce a global zero-point offset to use the parallaxes (suggested -0.04 mas)
- **You cannot correct the regional features:** if we could, we would already have corrected them. We have indications that these zero points may be present, but no more.
- For most of the sky assume an additional systematic error of 0.3 mas; your derived standard errors for anything cannot go below this value
 $\varpi \pm \sigma_{\varpi} \text{ (random)} \pm 0.3 \text{ mas (syst.)}$
- For a few smaller regions be aware that the systematics might reach 1 mas

More specifically: treat separately random error and bias, but if you must combine them, a **worst case** formula can be as follows

- **For individual parallaxes:** to be on the safe side add 0.3 mas to the standard uncertainty

$$\sigma_{\text{Total}} = \sqrt{\sigma_{\text{Std}}^2 + 0.3^2}$$

- **When averaging parallaxes for groups of stars:** the random error will decrease as \sqrt{N} but the systematic error (0.3 mas) will not decrease

$$\sigma_{\text{final}} = \sqrt{\sigma_{\text{averageStd}}^2 + 0.3^2}$$

where $\sigma_{\text{averageStd}}$ decrease is the formal standard deviation of the average, computed in the usual way from the sigmas of the individual values in the average (giving essentially the \sqrt{N} reduction).

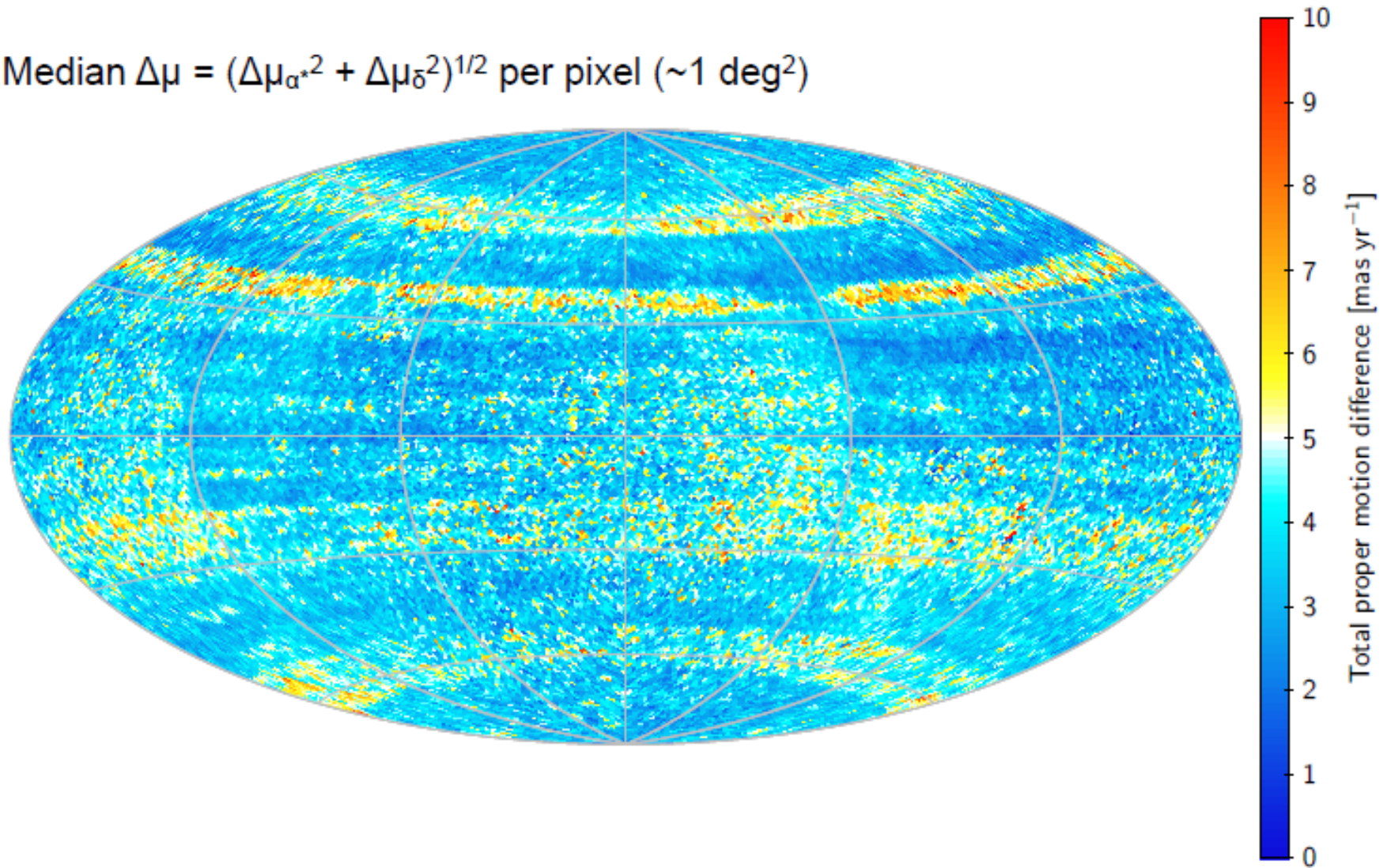
- **Don't try to get a “zonal correction” from previous figures, it's too risky**

For DR1 proper motions and positions:

- In this case Gaia data is the best available, by far.
- We do not have means to do a check as precise as the one done for parallaxes, but there are no indications of any significant bias
- For positions remember that for comparison purposes you will likely have to convert them to another epoch. You should propagate the errors accordingly.

Comparison with Tycho-2 shows that catalogue's systematics (not Gaia's)

Median $\Delta\mu = (\Delta\mu_\alpha^2 + \Delta\mu_\delta^2)^{1/2}$ per pixel ($\sim 1 \text{ deg}^2$)



Errors 2: random errors

Random error:

your measurements are randomly distributed around the true value

- Each measurement in a catalogue comes with a formal error
- Random errors usually are quasi-normal.
- The formal error is meant to represent the variance of a normal distribution around the true value

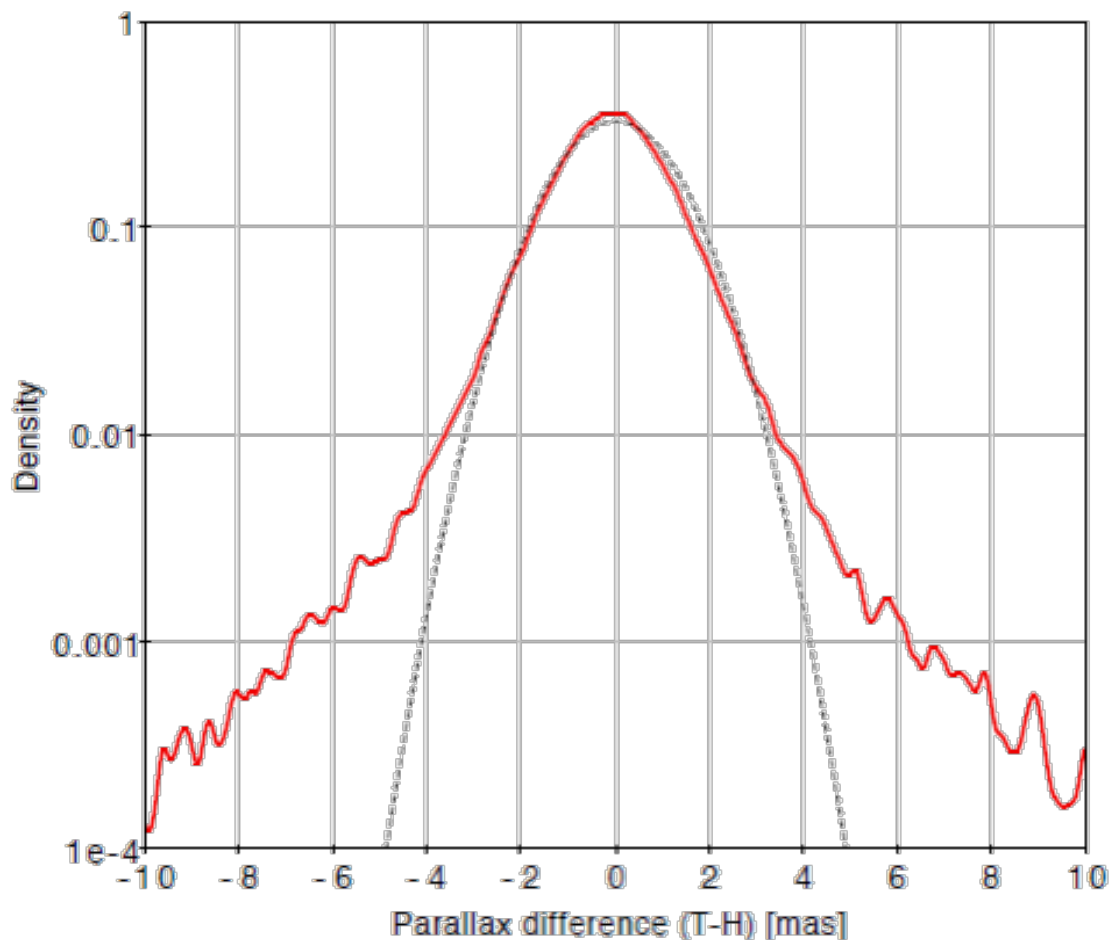
Example:

- Published formal errors for Gaia DR1 may be slightly overestimated
- However, in most scientific data sets they are underestimated

Warning 1: Outliers

comparison with Hipparcos shows deviation from normality beyond ~ 2 ?

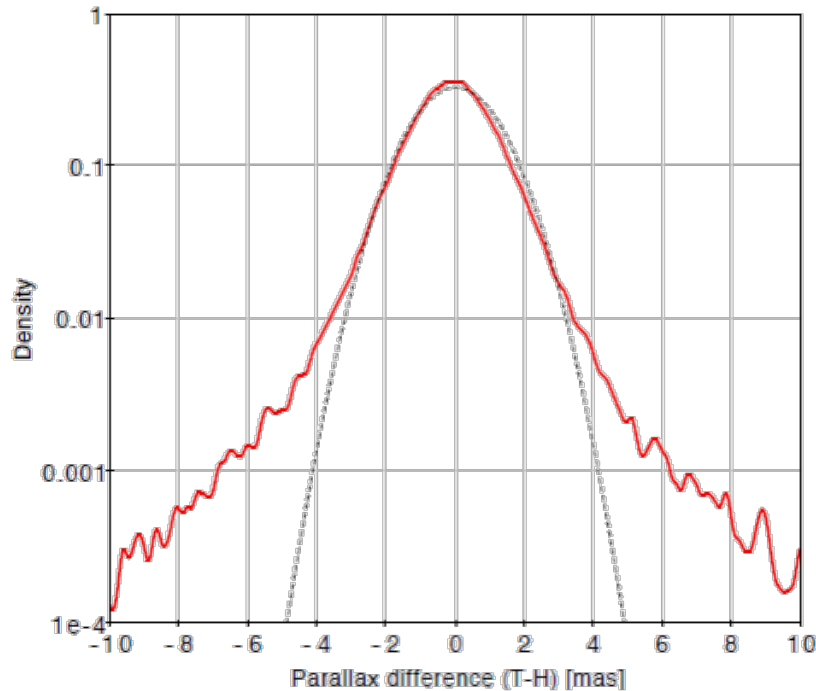
$\text{med}(\Delta\varpi) = -0.086 \text{ mas}, \text{RSE} = 1.22 \text{ mas}$



To take into account for outlier analysis

Warning 1: non-Gaussianity; outliers

$\text{med}(\Delta\varpi) = -0.086 \text{ mas}$, $\text{RSE} = 1.22 \text{ mas}$

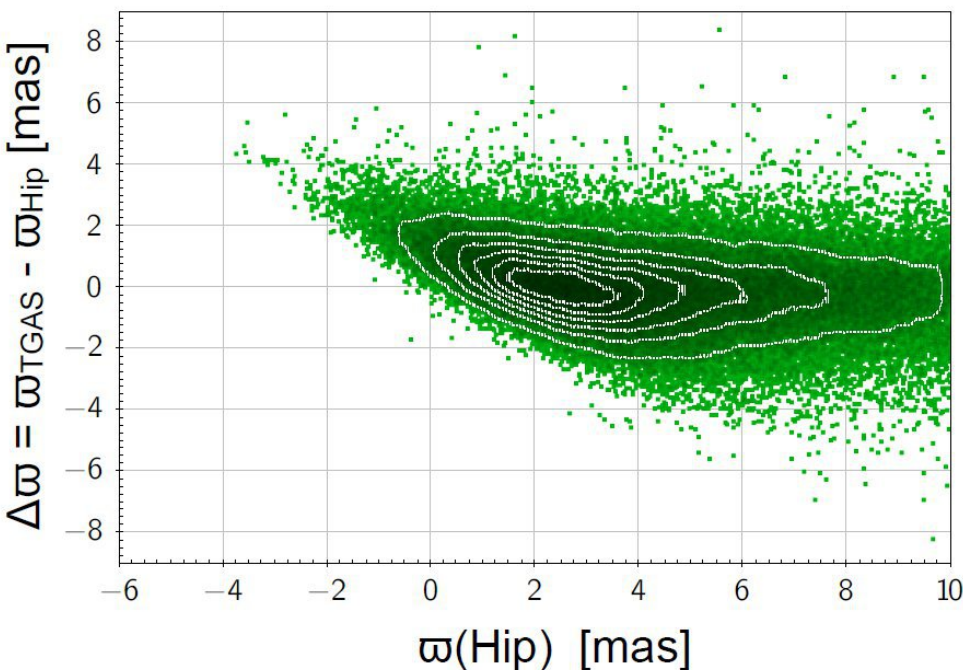


- **Comparison TGAS vs Hipparcos:**
deviation from normality beyond ~ 2.5 ?
- **TGAS negative parallaxes:**
a long negative tail is apparent
- **How to take into account:**
always do an outlier analysis
(if possible ...)

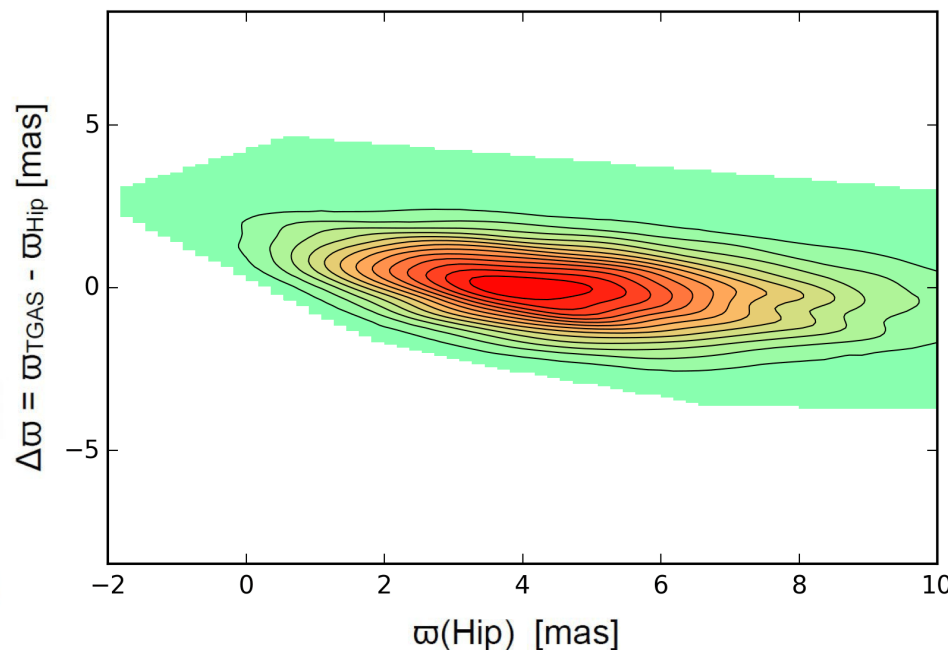
Warning 2: when comparing with other sources of trigonometric parallaxes take into account the properties of the error distributions

TGAS vs Hipparcos

Observations



Simulations

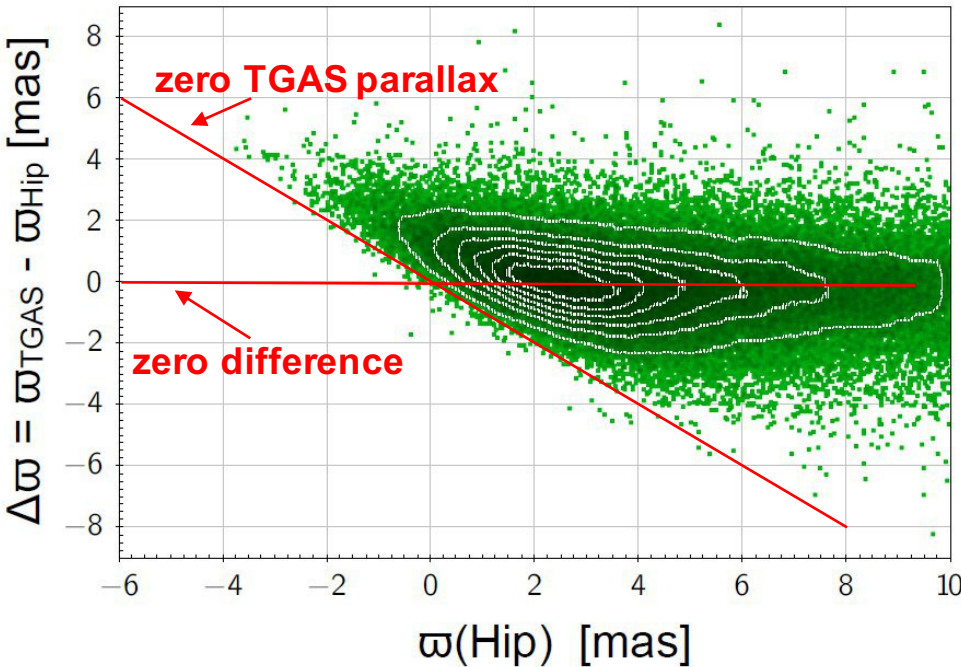


The “slope” at small parallaxes is not a bias in either TGAS or HIP, simply due to the different size of the errors in the two catalogues!

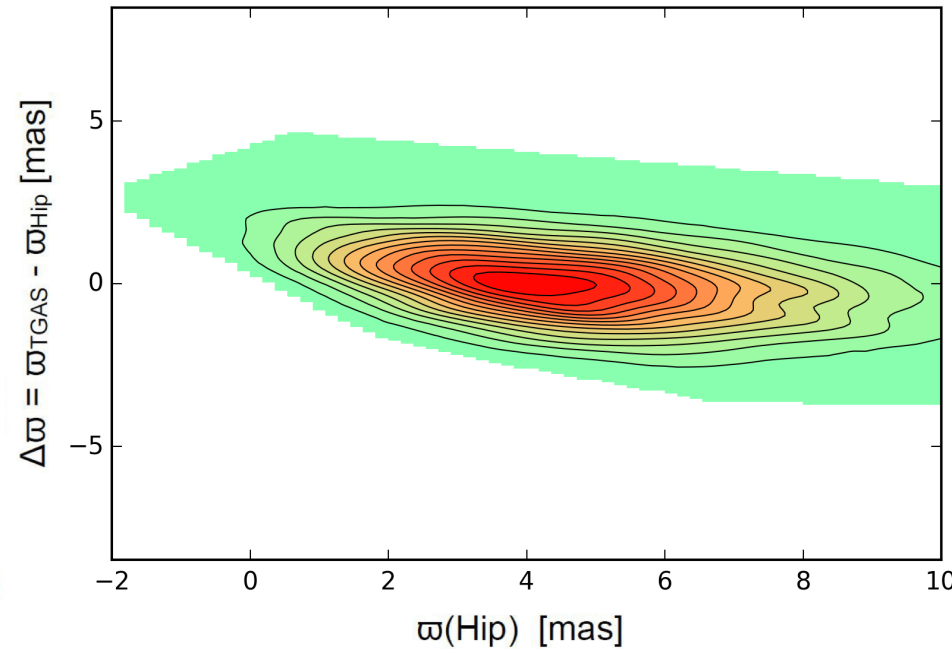
Warning 2: when comparing with other sources of trigonometric parallaxes take into account the properties of the error distributions

TGAS vs Hipparcos

Observations



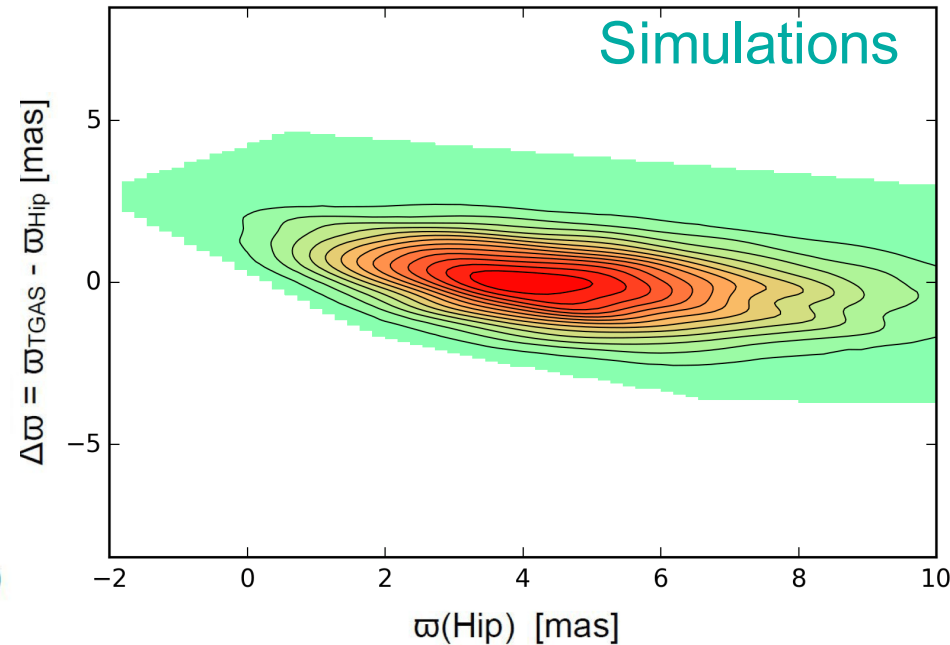
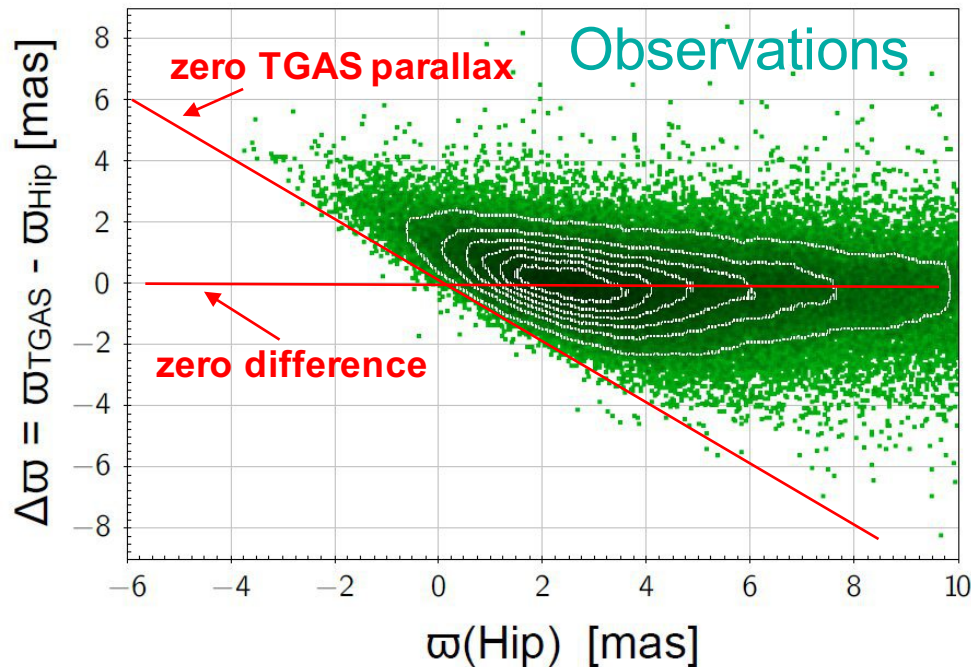
Simulations



The “slope” at small parallaxes is not a bias in either TGAS or HIP, simply due to the different size of the errors in the two catalogues!

Warning 2: spurious biases

Example 1: Comparison TGAS vs Hipparcos

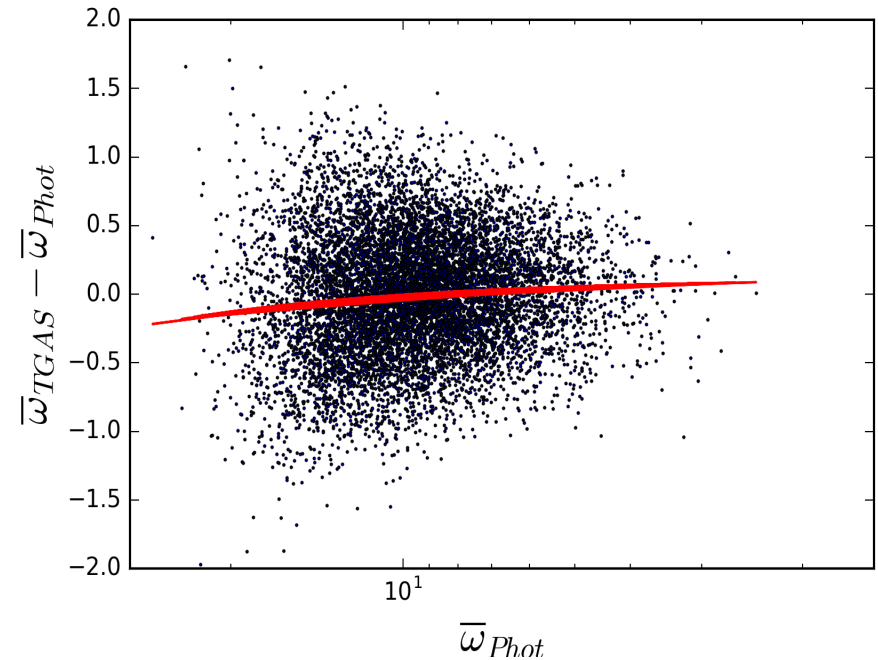
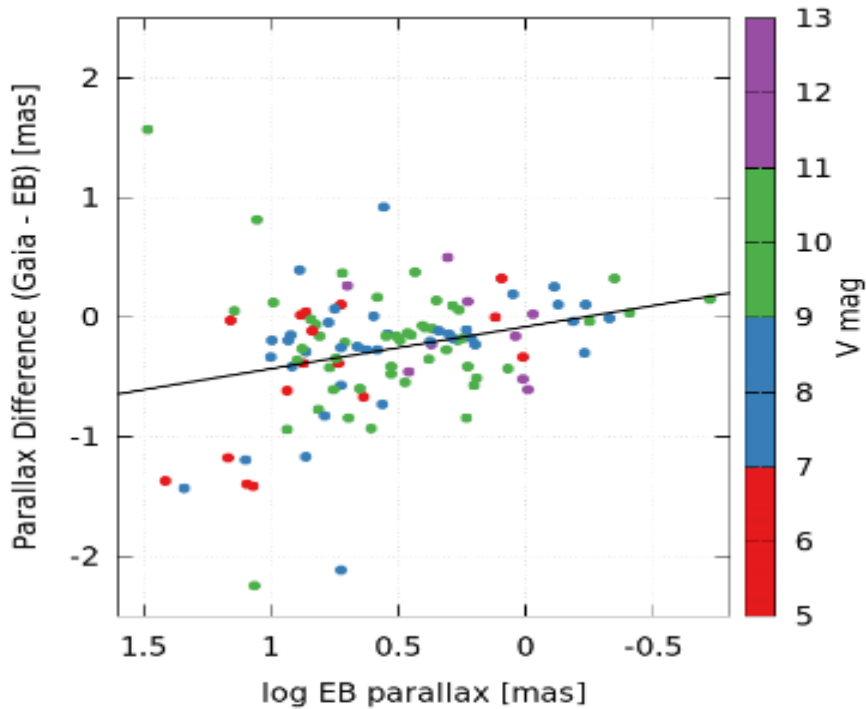


- The “slope” at small parallaxes is not a bias in either TGAS or HIP:
It is simply due to the different size of the errors in the two catalogues!
- How to take into account:
always consider the **widths** of the error distributions

Example 2: Eclipsing binaries parallaxes vs TGAS

arXiv:1609.05390v3

Simulation



The overall “slope” is due to the different **shapes** of the error distributions in parallax
(log-normal for photometric, normal for trigonometric)

Errors 3: correlations

Correlation:

the measurements of several quantities are not independent from each other

- Whenever you take linear combinations of such quantities, the correlations have to be taken into account in the error calculus (and even more so for non-linear functions)
- Example:
 - The errors in the five astrometric parameters for each source in Gaia DR1 are not independent of each other
 - Therefore the ten correlations between these parameters are provided (correlation matrix)
 - Use cases:
Galactic proper-motion components, positions after epoch transformation, ...
- How to: for recipe(s) see the omitted pages on the presentations folder

Errors 3: correlations

Correlation:

the measurements of several quantities are not independent from each other.

- Whenever you take linear combinations of such quantities, the correlations have to be taken into account in the error calculus (and even more so for non-linear functions !)

Variance of a sum: (x_1+x_2)

$$\begin{aligned}\sigma^2(x_1+x_2) &= \sigma^2(x_1) + \sigma^2(x_2) + 2 \operatorname{cov}(x_1,x_2) \\ &= \sigma^2(x_1) + \sigma^2(x_2) + 2 \sigma(x_1) \sigma(x_2) \operatorname{corr}(x_1,x_2)\end{aligned}$$

Variance of any linear combination of two measured quantities, x_1 and x_2 : $(ax_1 + bx_2)$

$$\begin{aligned}\sigma^2 &= a^2 \sigma^2(x_1) + b^2 \sigma^2(x_2) + 2ab \operatorname{cov}(x_1,x_2) \\ &= a^2 \sigma^2(x_1) + b^2 \sigma^2(x_2) + 2ab \sigma(x_1) \sigma(x_2) \operatorname{corr}(x_1,x_2)\end{aligned}$$

Generally, for a whole set of linear combinations y of several correlated random variables x :

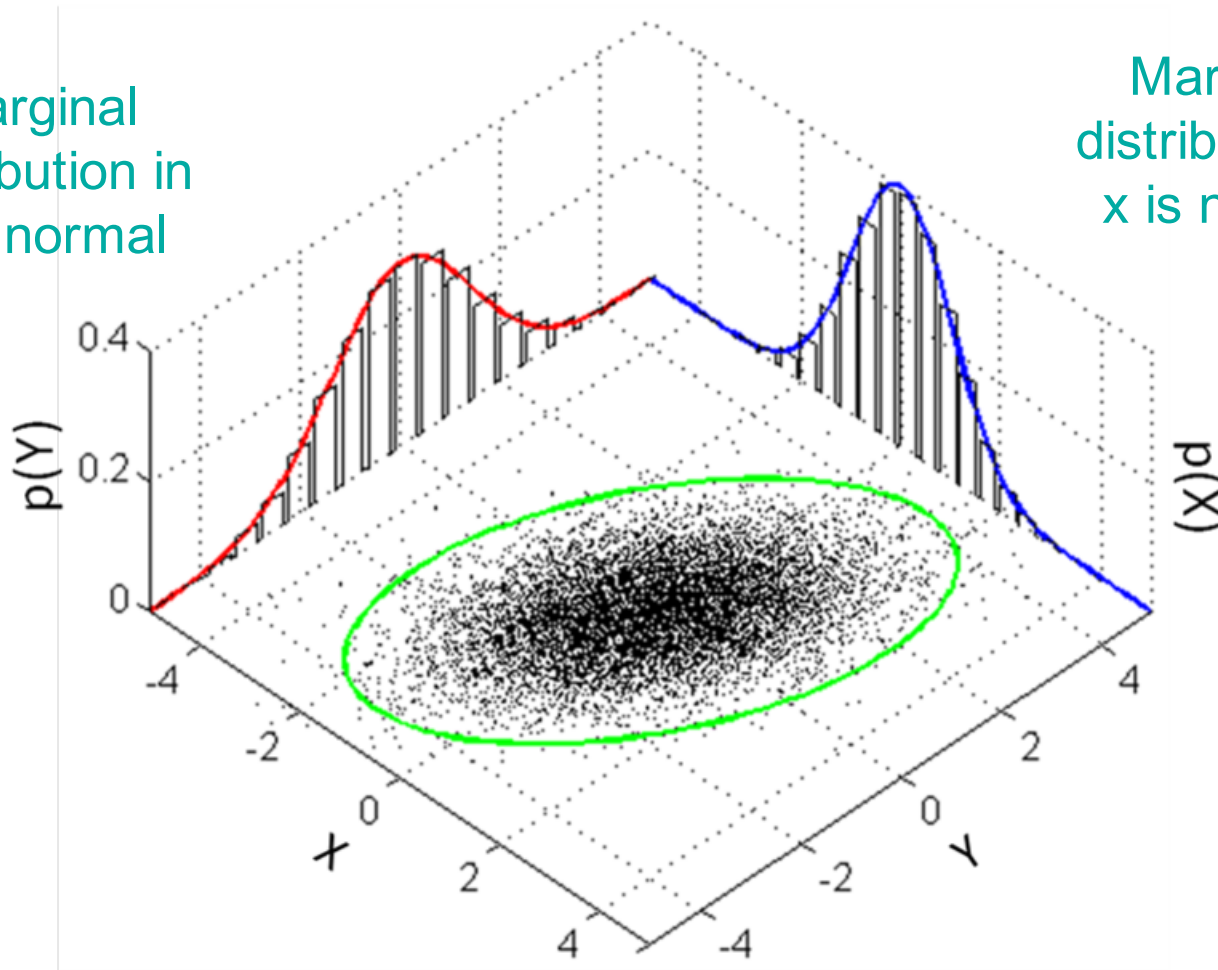
If $y = A'x$, then: $\operatorname{Cov}(y) = A' \operatorname{Cov}(x) A = A' \operatorname{Sigma}(x) \operatorname{Corr}(x) \operatorname{Sigma}'(x) A$

where Cov and Corr indicate covariance and correlation matrices, $\operatorname{Sigma}(x)$ is a diagonal matrix having the sigmas of the components of x as elements, and A' is the relation matrix. In the example above, for just two x and one y , the matrix A' is simply the row vector (a,b) .

Example of two correlated parameters

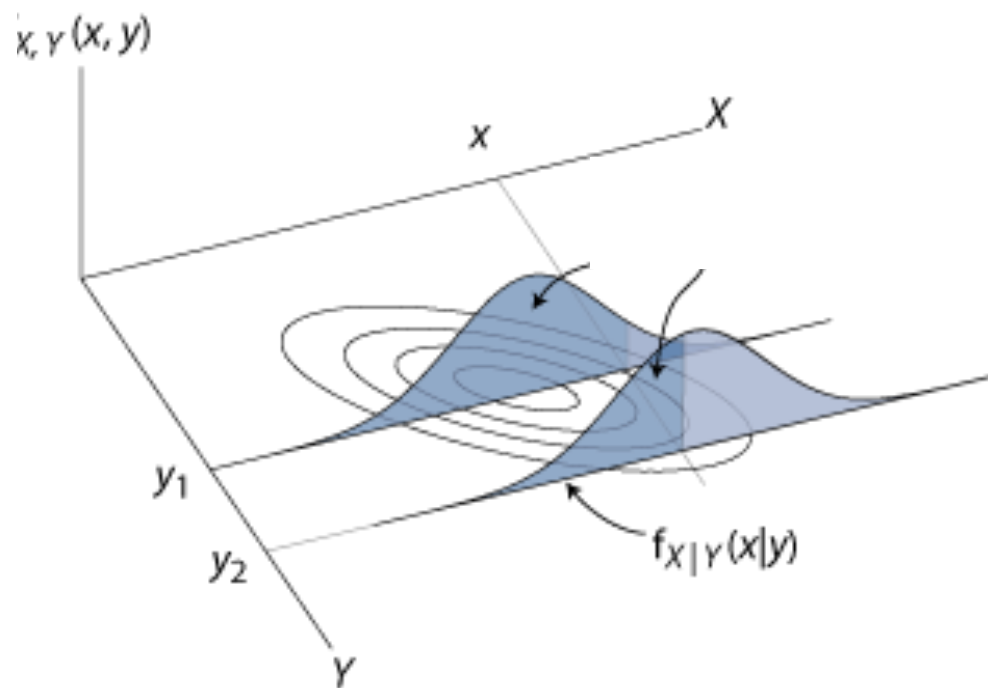
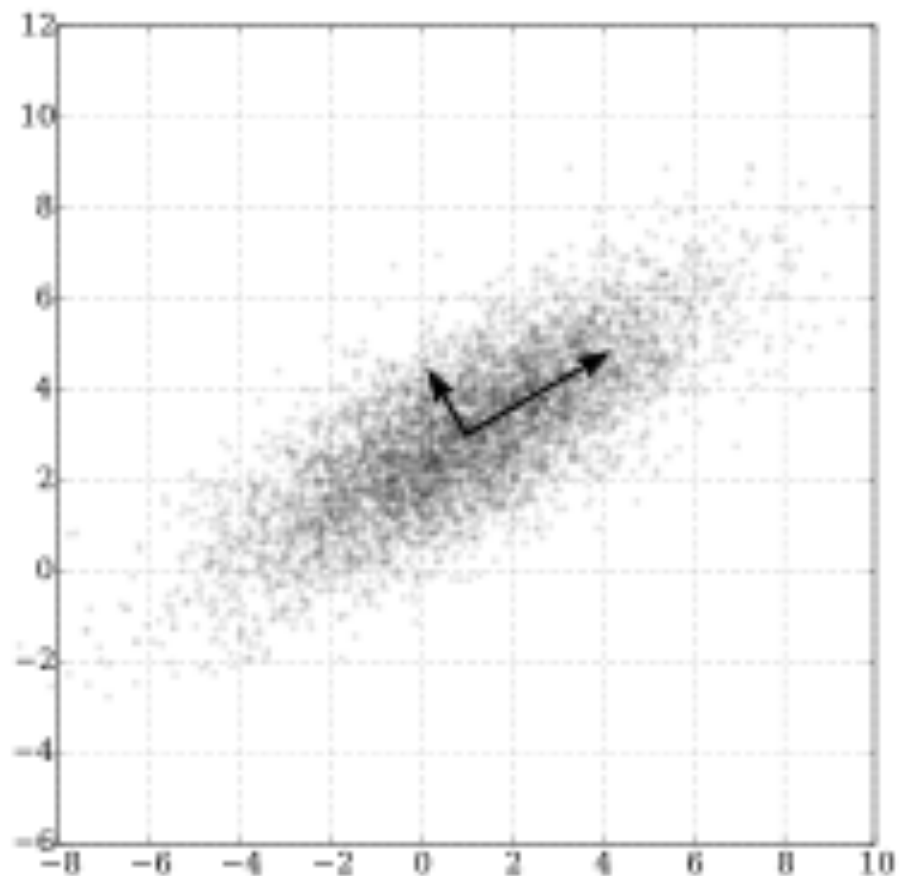
Marginal distribution in y is normal

Marginal distribution in x is normal



By Bscan - Own work, CC0, <https://commons.wikimedia.org/w/index.php?curid=25235145>

Beware when using these quantities together



Examples of problematic use:

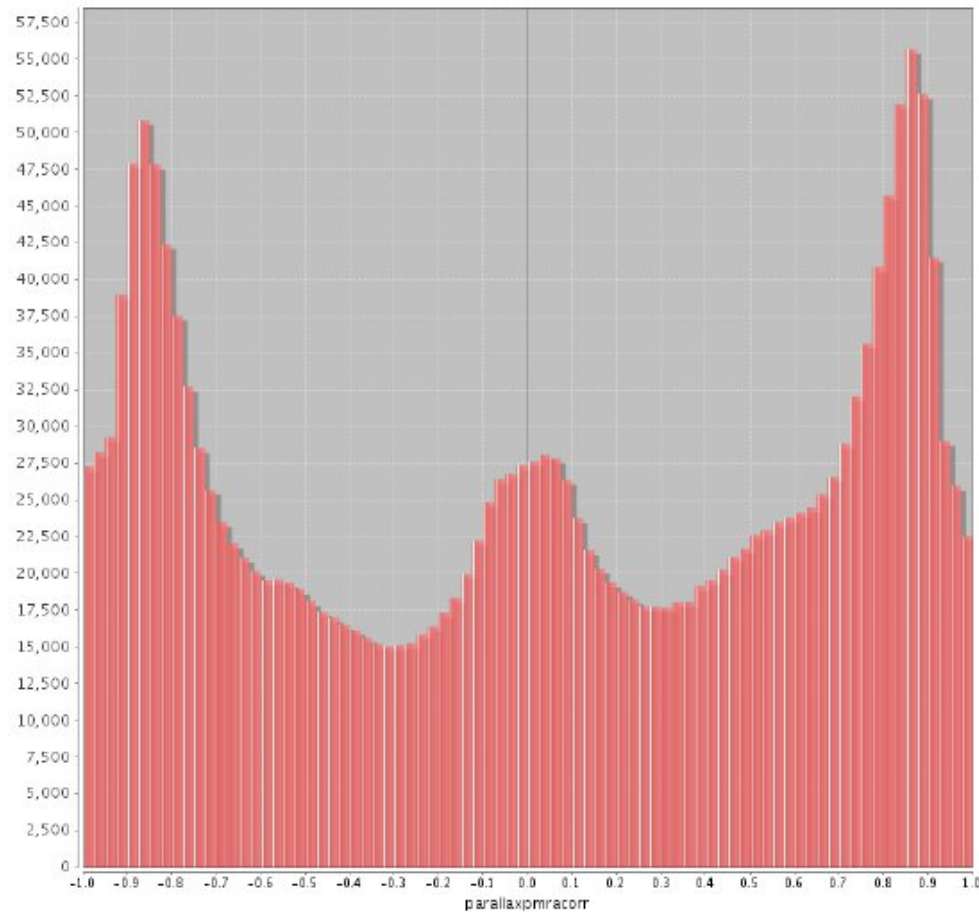
- Simple epoch propagation (!) pos&pm
- Calculation of proper directions pos&pm¶llax
- Proper motion in a given direction on the sky (other than north-south or east-west) proper-motion components
- Proper motion components in galactic or ecliptic coordinates proper-motion components

- More complex, non-linear example:

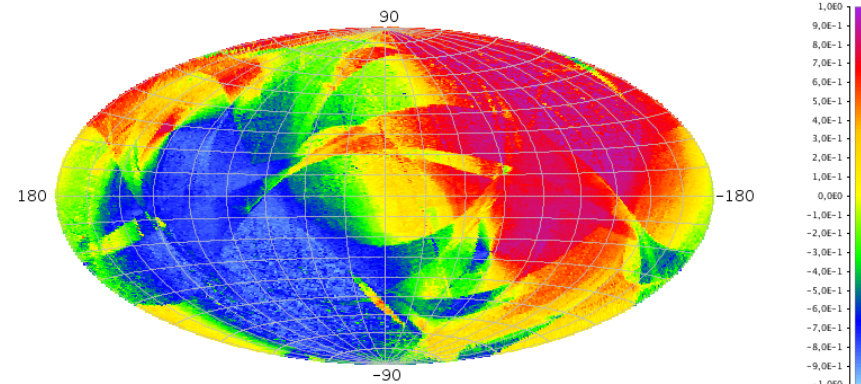
Calculating the transversal velocities of a set of stars

- The resulting dispersion of velocities is influenced by the errors in parallax and in proper motion; thus 3-dimensional case.
- Its determination can not be done using the parallax and proper motion errors separately; the correlations have to be taken into account
- But this time it's non-linear! The error distribution will no longer be Gaussian.
- The A matrix of the previous page will become the Jacobian matrix of the local derivatives of the transversal velocity wrt parallax and pm components

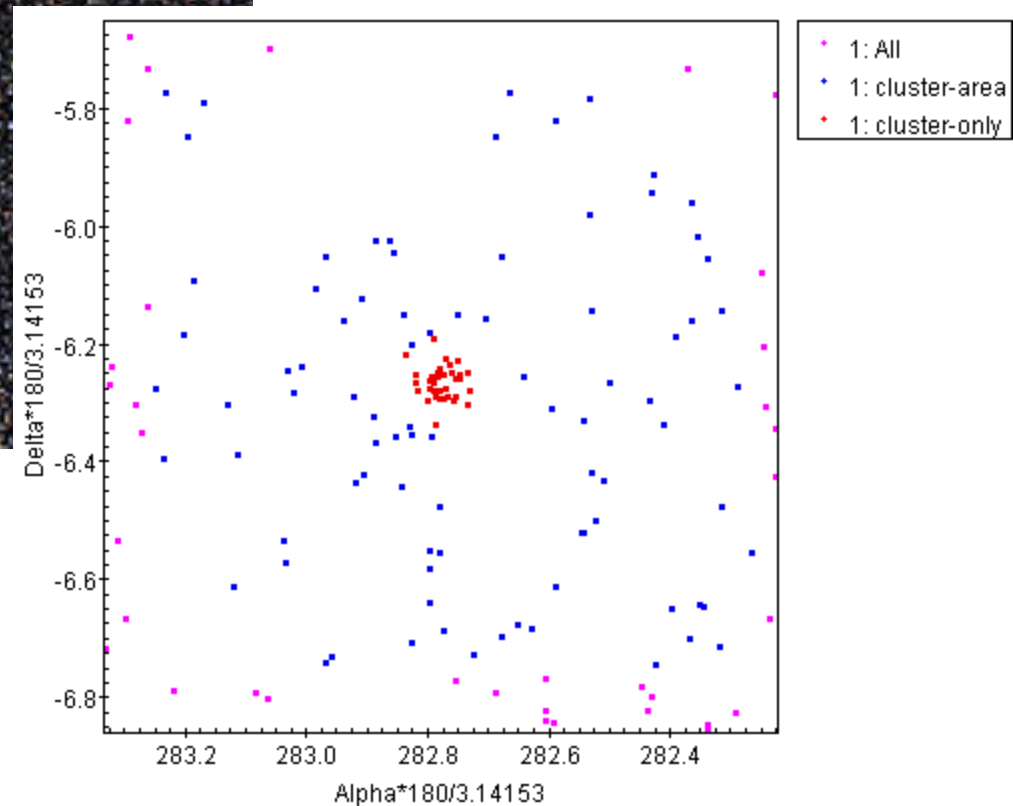
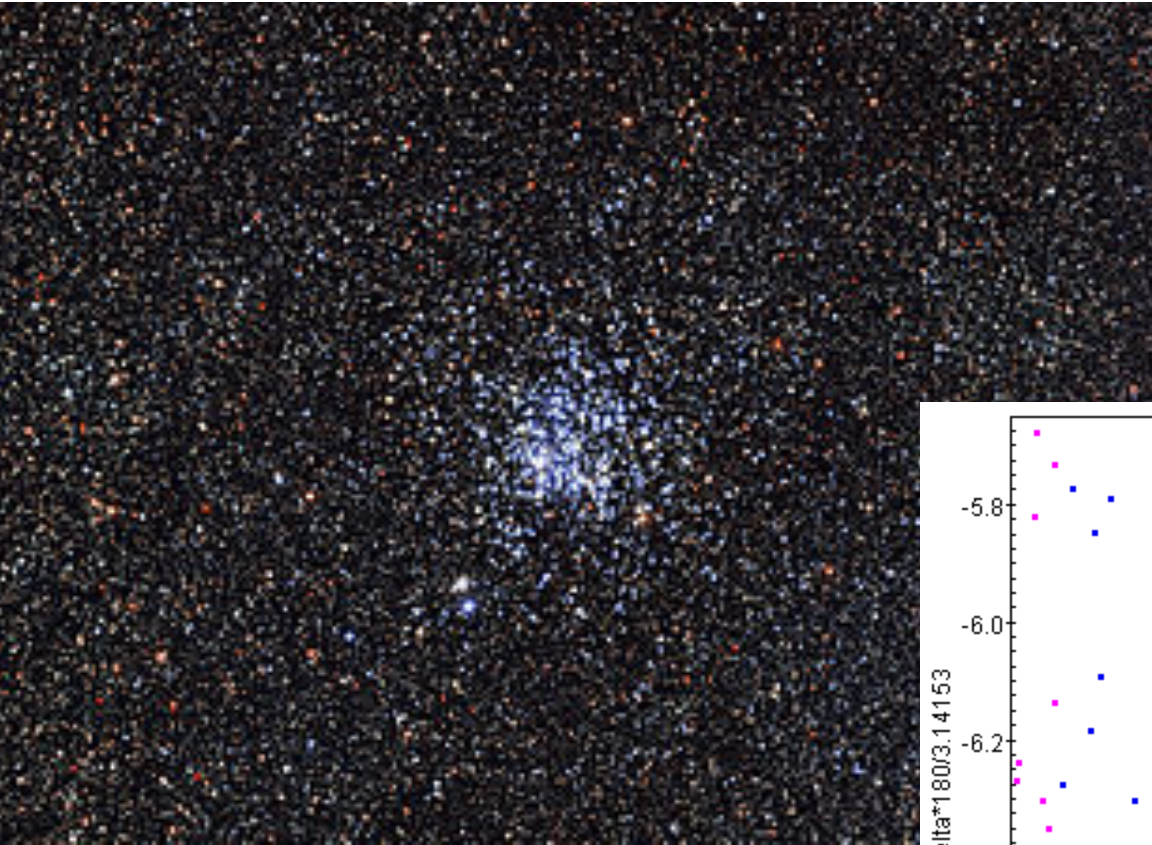
Beware: large and unevenly distributed correlations in DR1; example: PmRA-vs.-Parallax correlation



HealpixMapMean parallax and pmra correlation in GAL coordinates (Value of objects). Objects: 2057050. Objects Out: 0



A really pretty example on correlations: M11



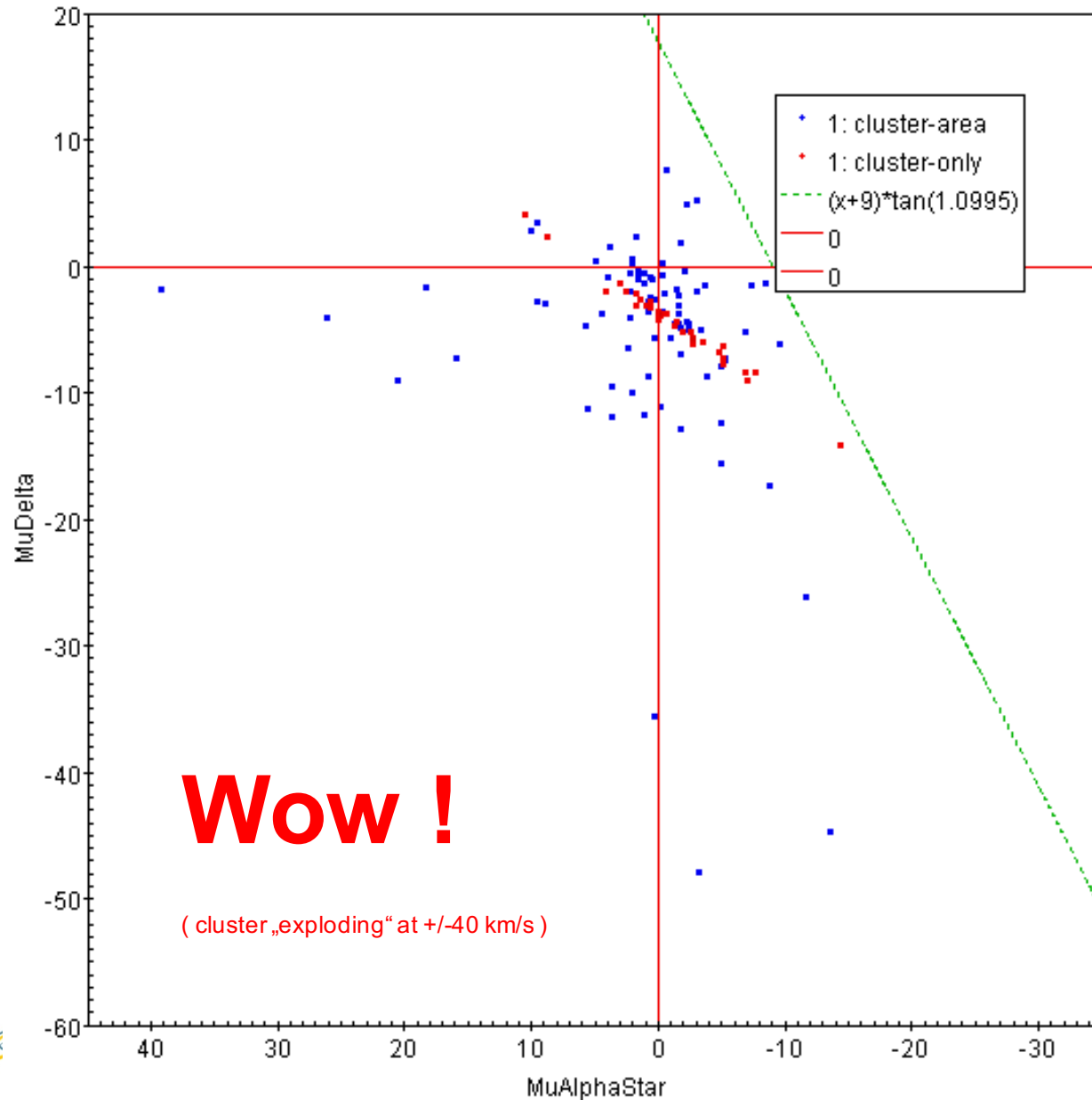
gaia



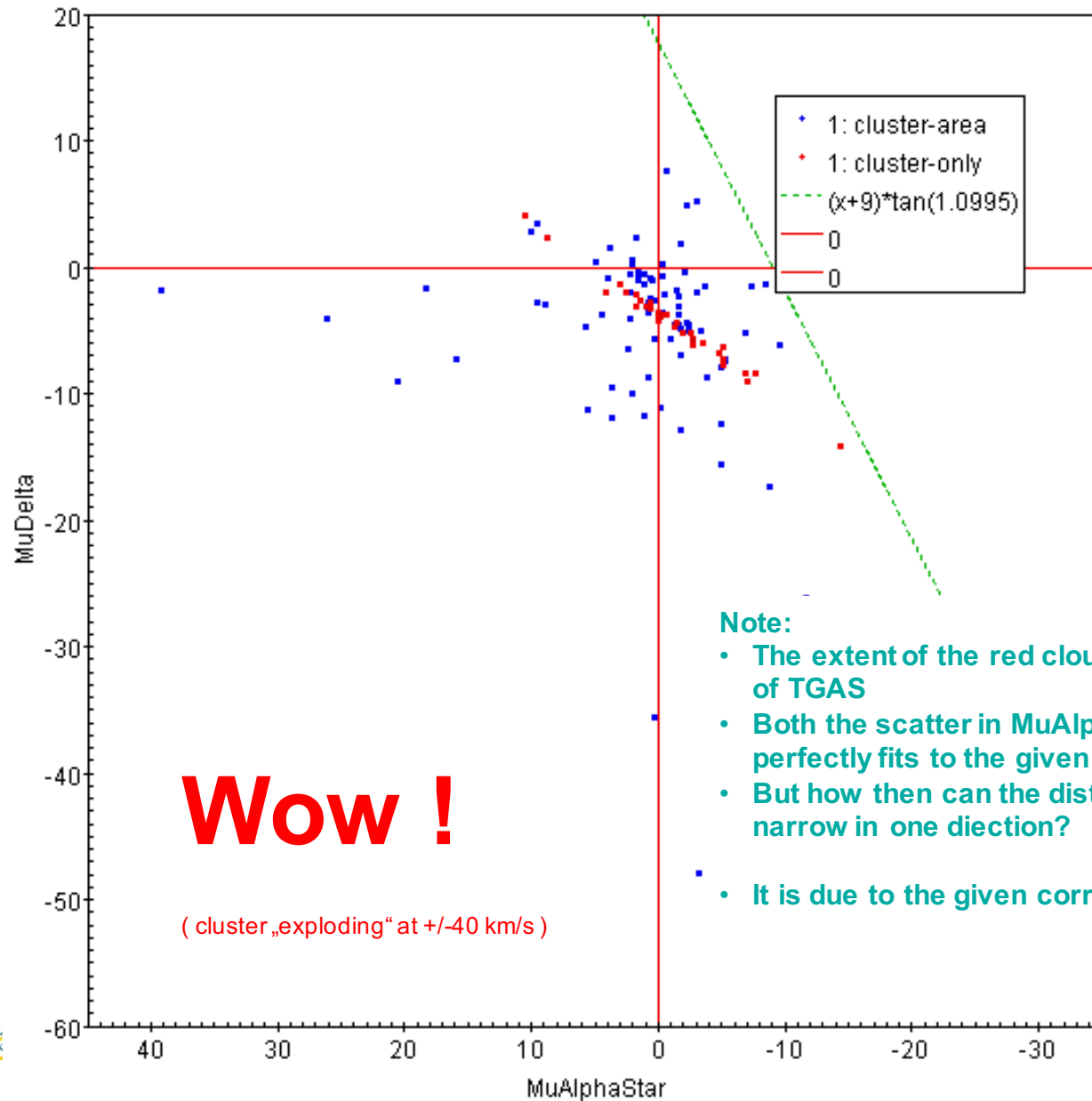
Gaia
DPAC
Data Processing & Analysis Consortium



M11; proper motions in the AGIS-01 solution



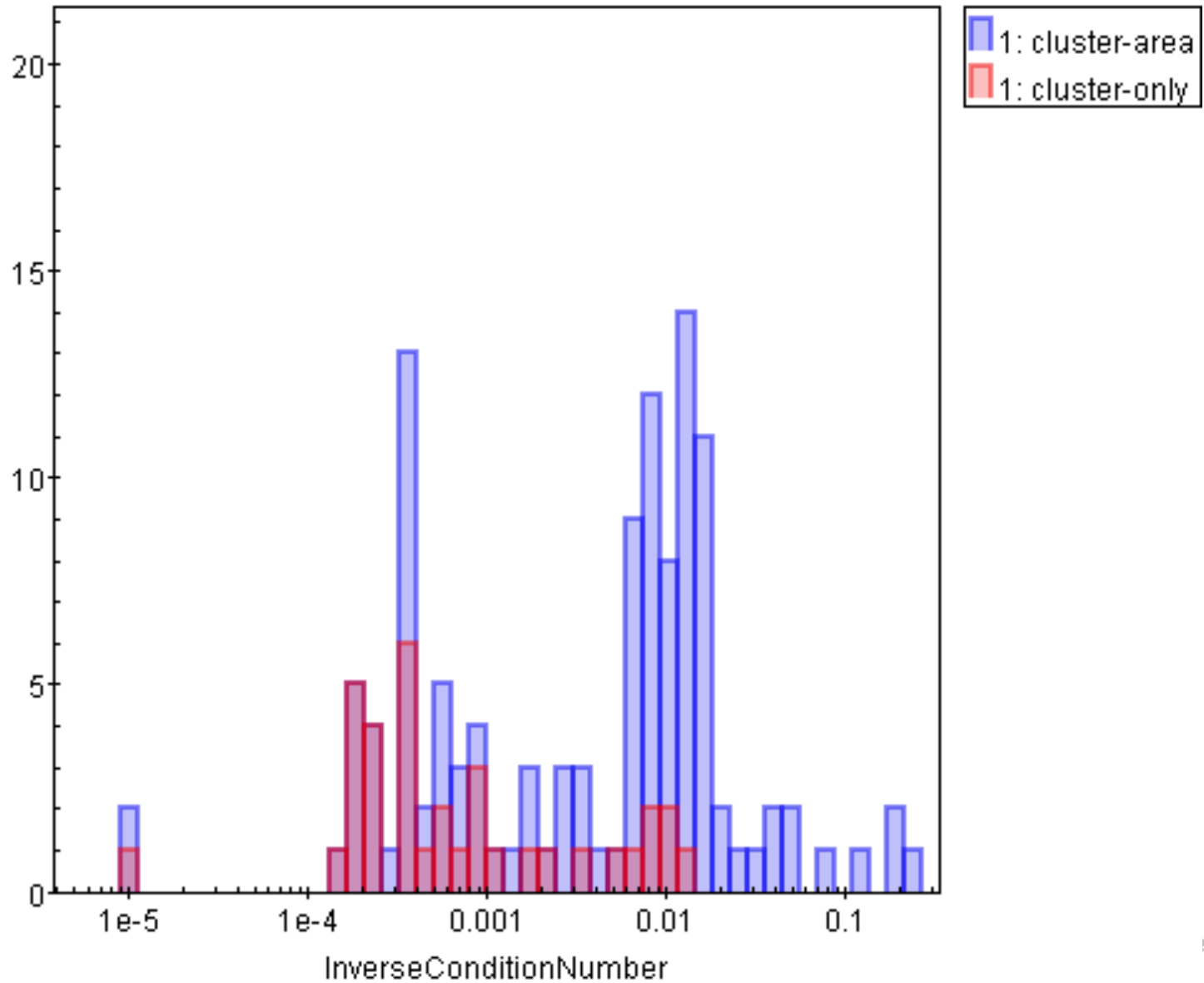
M11; proper motions in the AGIS-01 solution



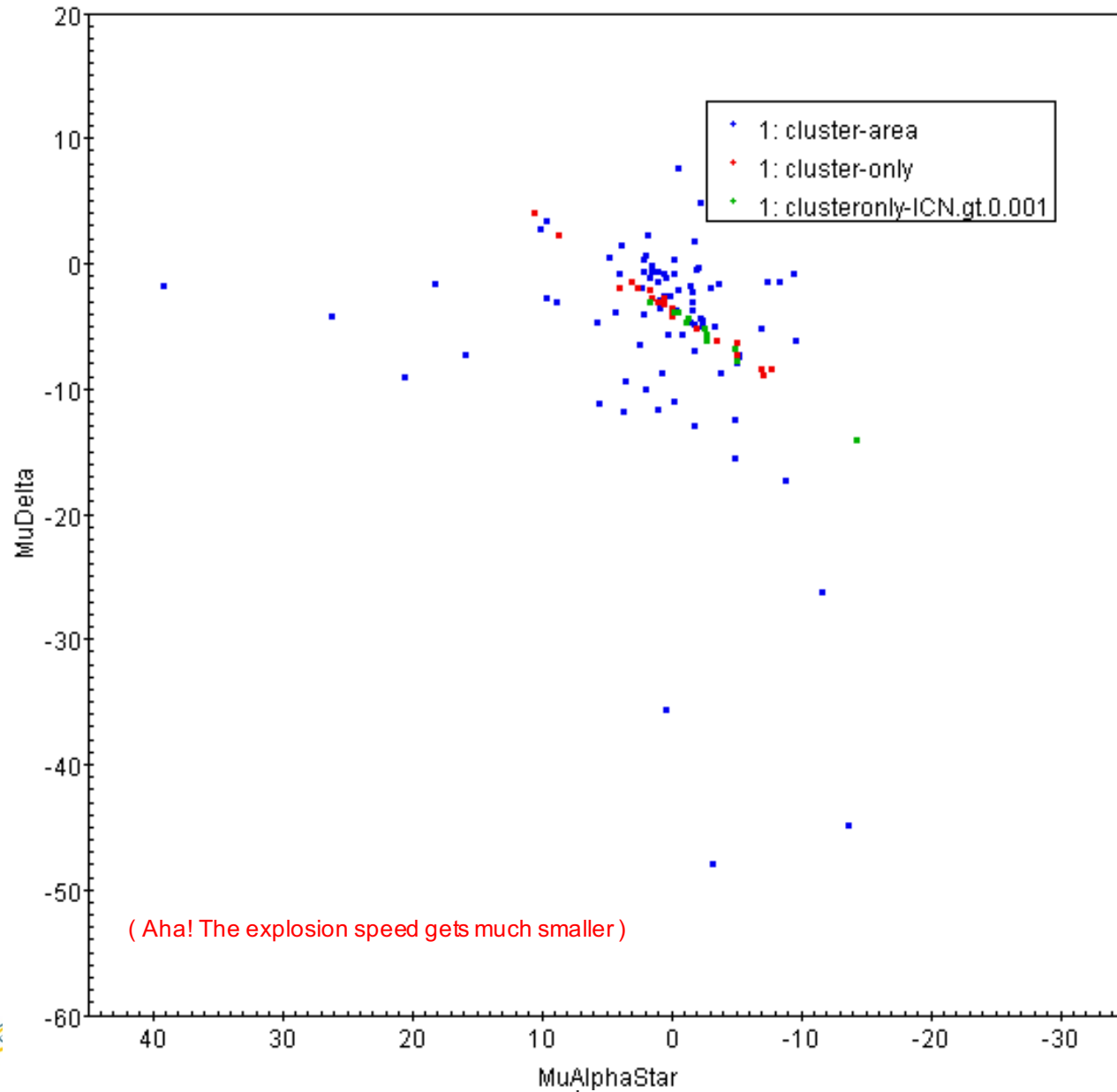
Wow !

(cluster „exploding“ at +/-40 km/s)

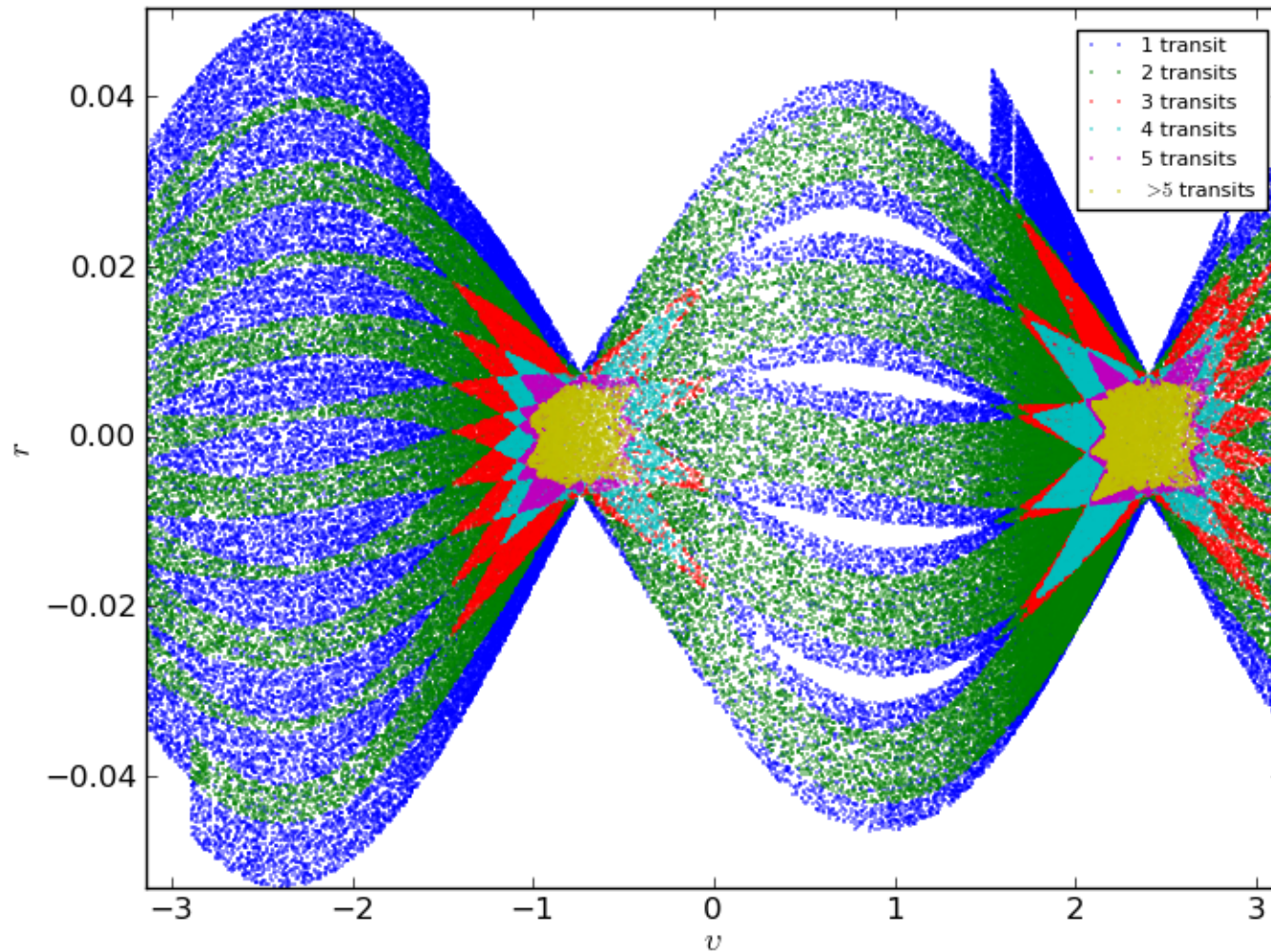
M11; scan coverage statistics



M11; selection of „better-observed“ stars



Just bad luck for poor M11:



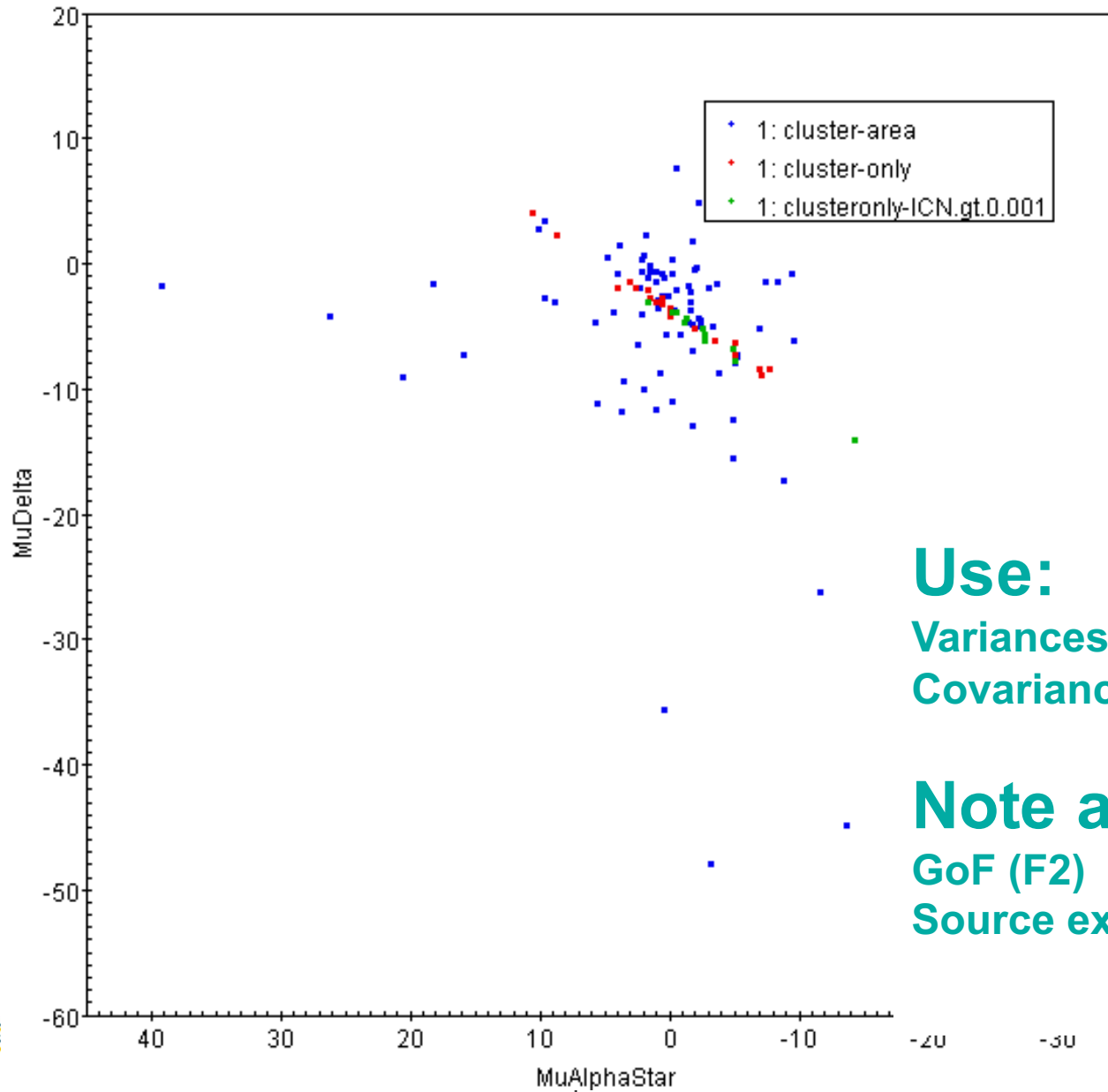
6 transits

all but one ...

slits

hickups

M11; lessons to be learned



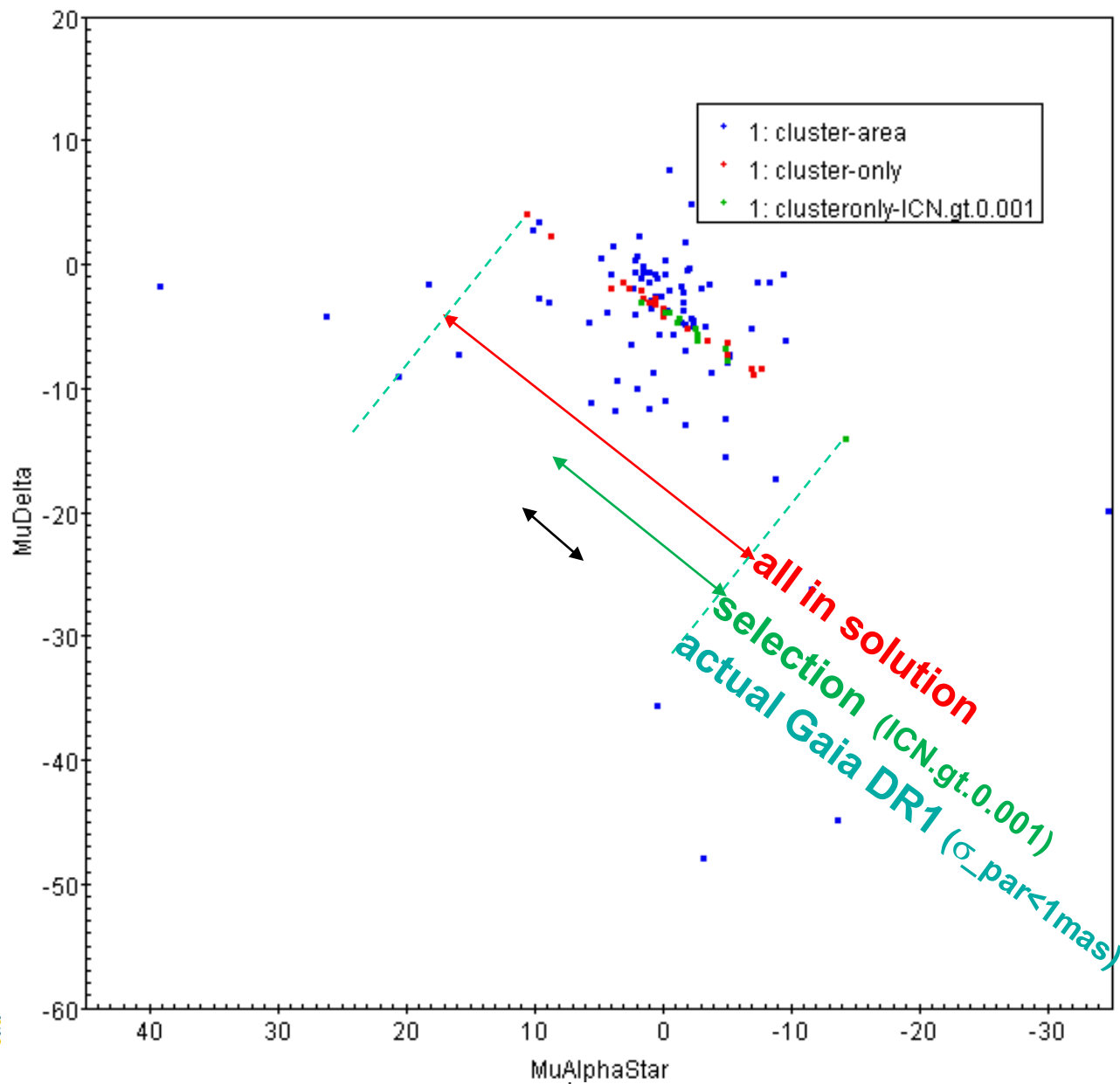
Use:

Variations/mean errors
Covariances/Correlations

Note and use:

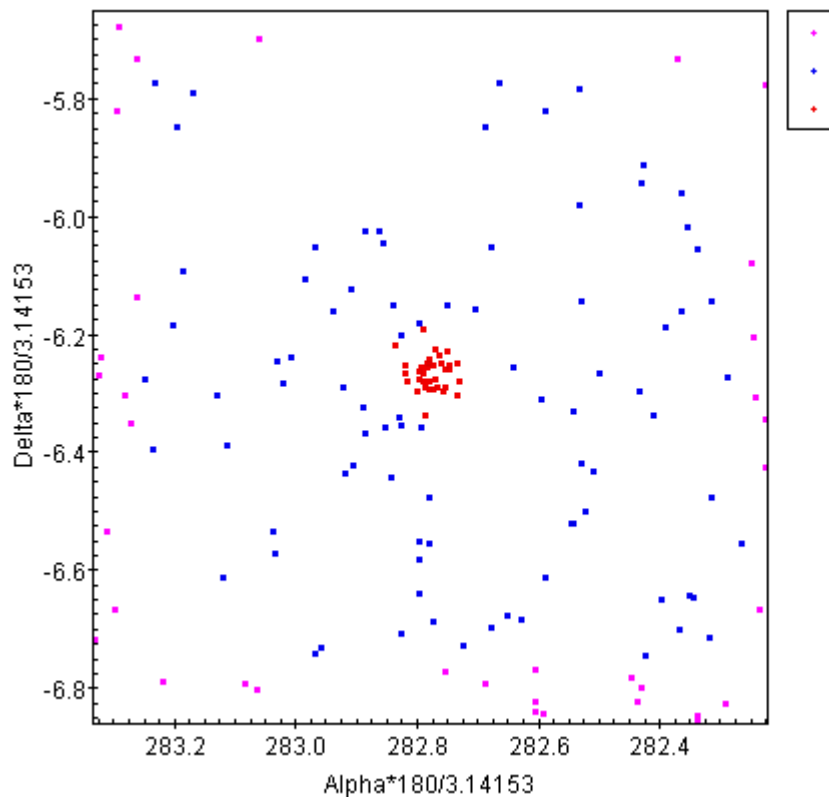
GoF (F2)
Source excess noise

M11; reasonable selection improves things

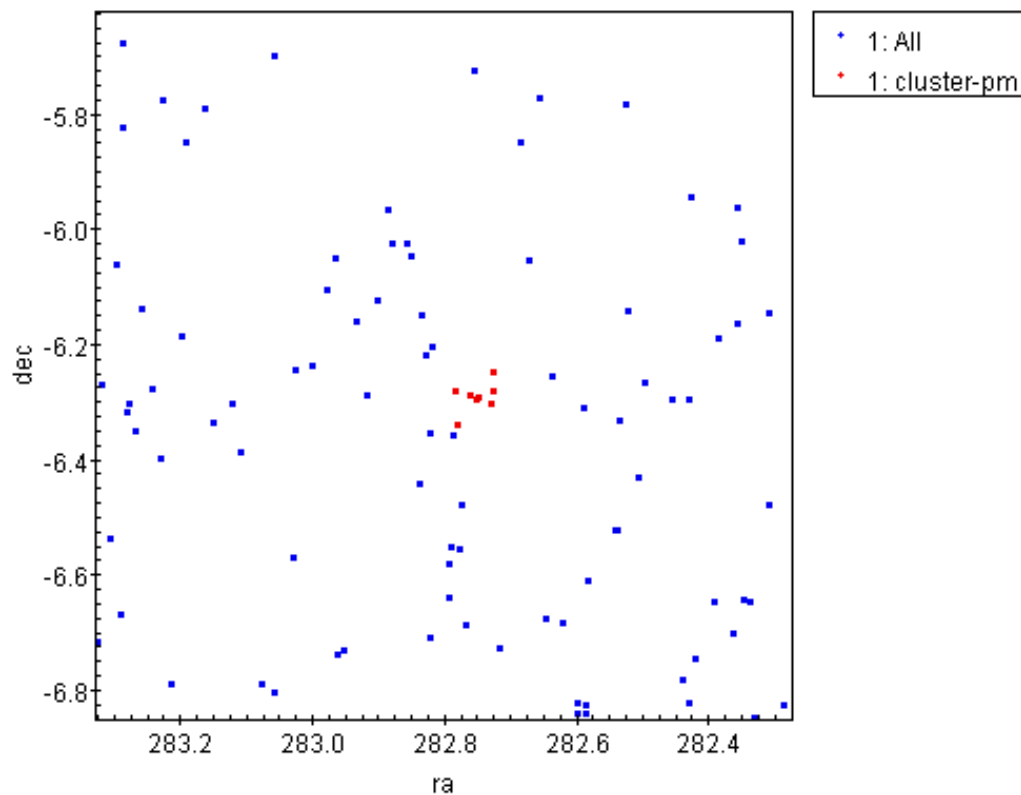


But there's always a price to be payed:

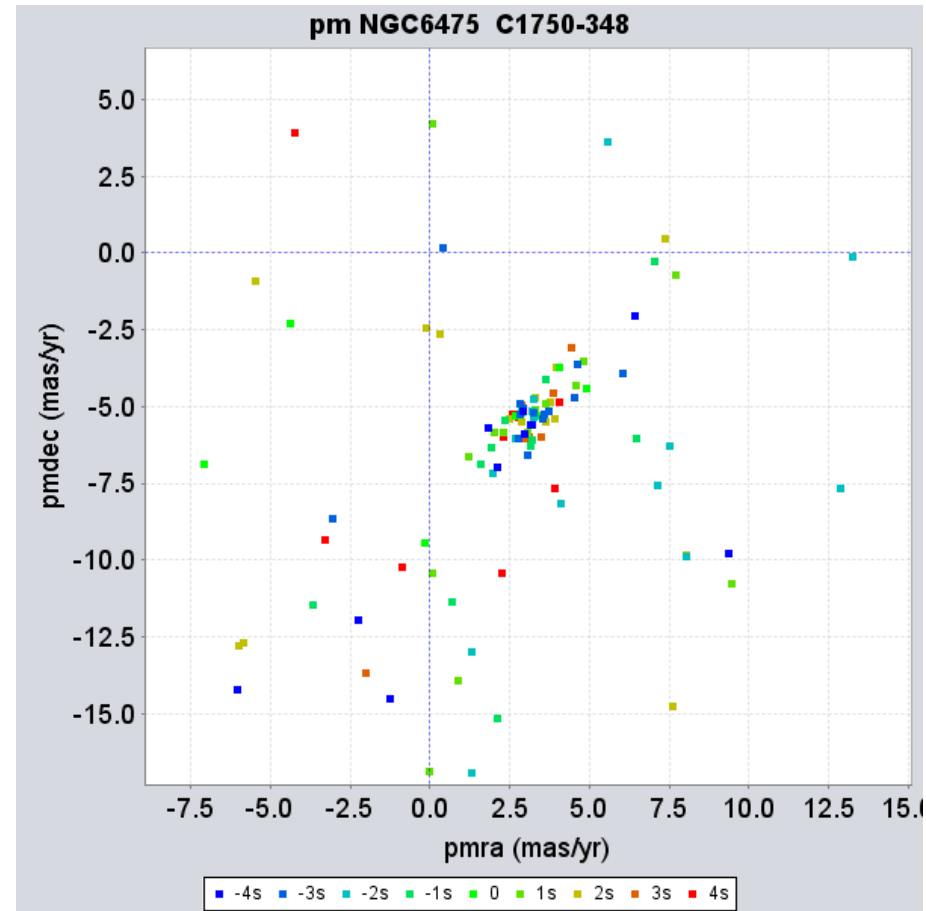
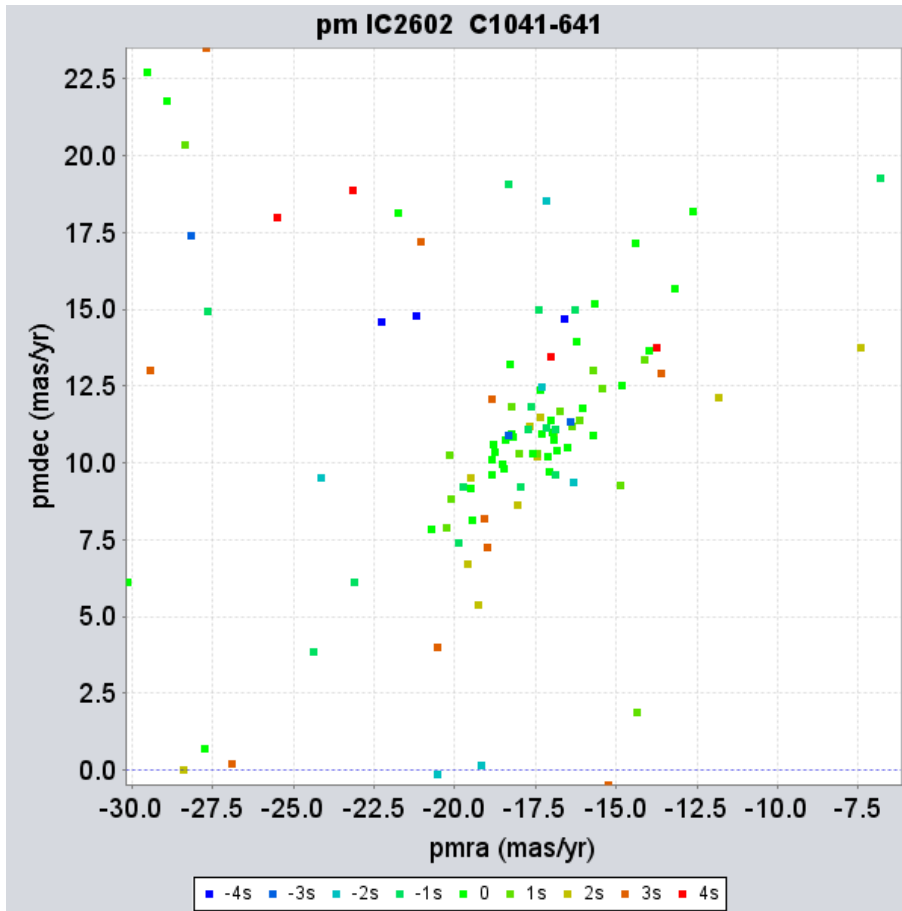
all in TGAS solution



actually in Gaia DR1



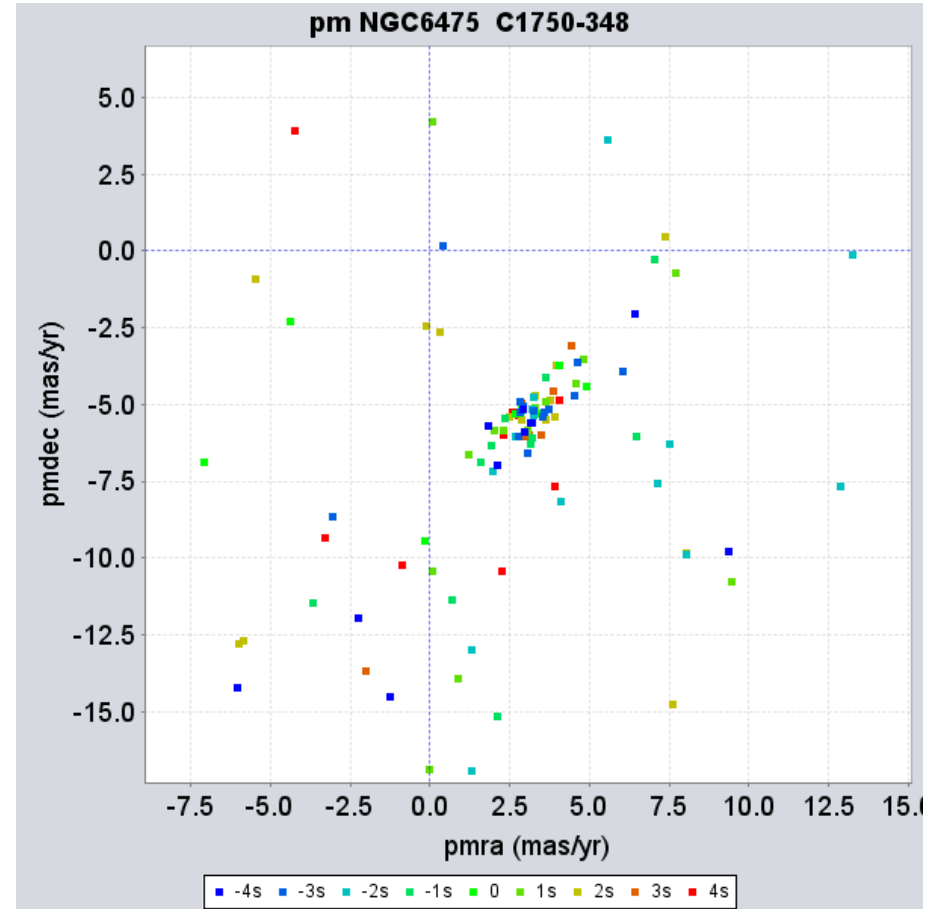
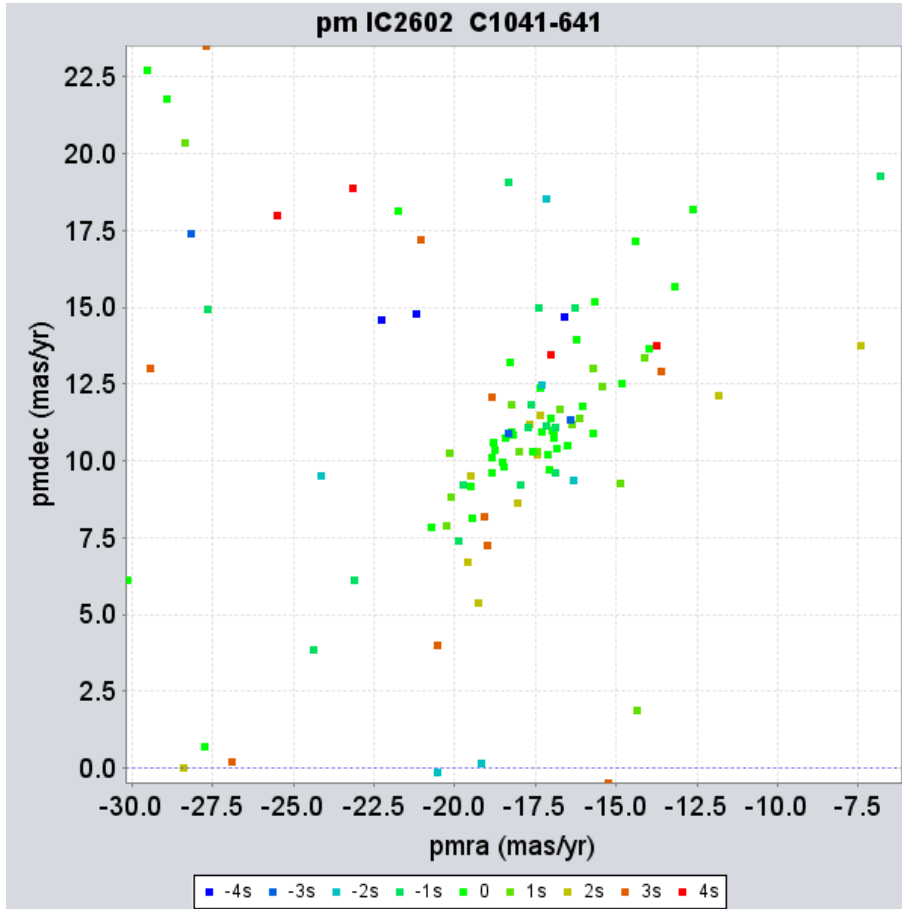
M11 is an extreme case, but ...



Two less extreme but still clearcut cases; using public DR1 data.

Note: the scales of the two figures are equal. NGC 6475 measured much more precisely.

Example: Star clusters seemingly „exploding“



Public DR1 data.

Note: the scales of the two figures are equal. NGC 6475 measured much more precisely.

Chapter 4: Transformations

Transformations:

when the quantity you want to study
is not the quantity you observe

Examples:

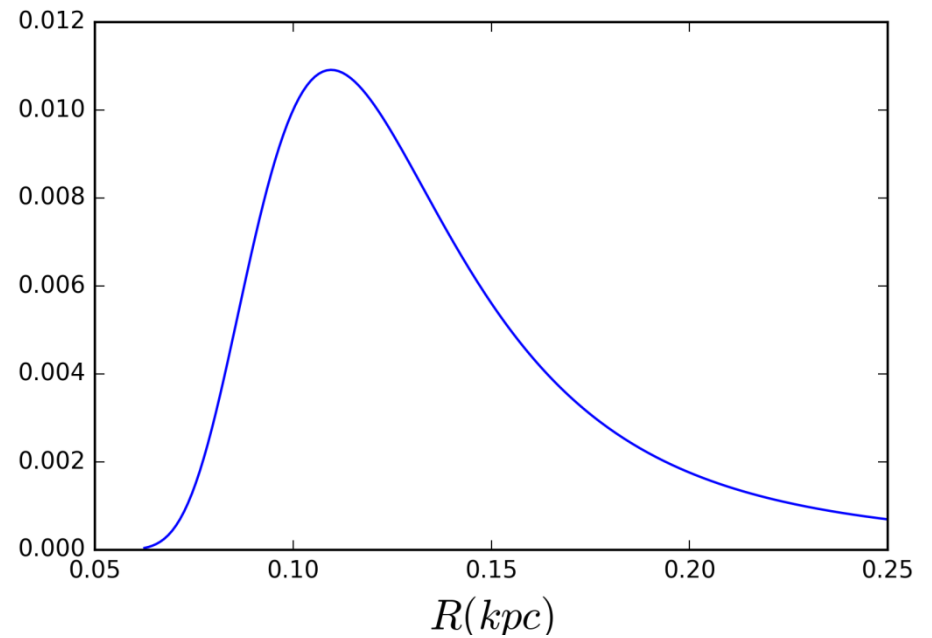
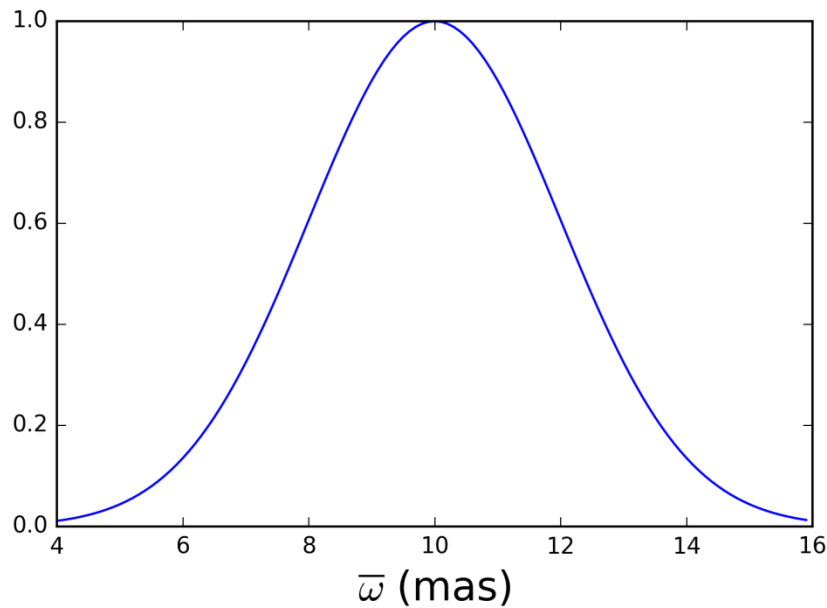
- Usually you want distances, not parallaxes
- Usually you want spatial velocities, not proper motions

Warning:

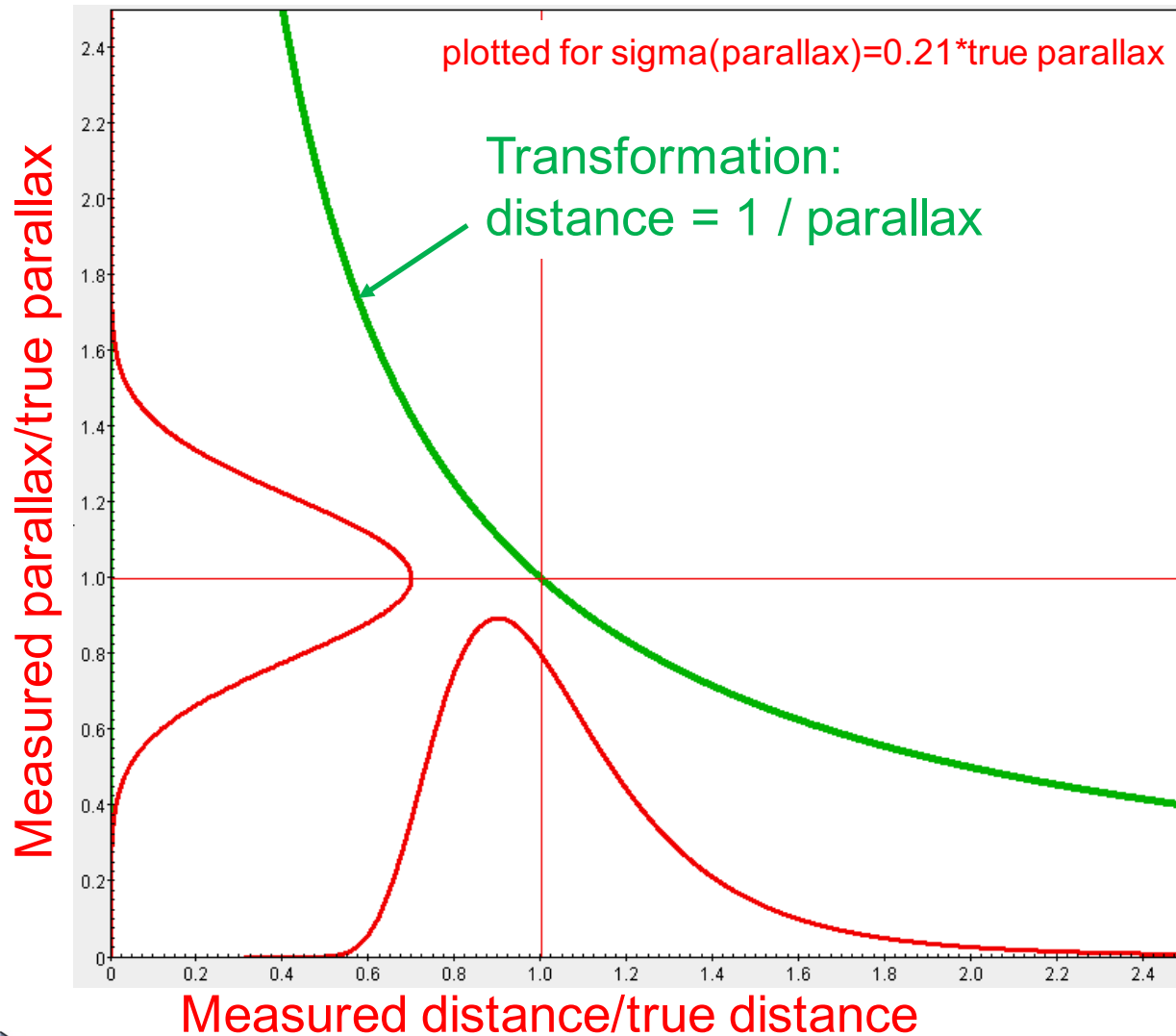
when using a transformed quantity the error distribution also is transformed

- This is especially crucial for the calculation of distances from parallaxes
- And even more so for the calculation of luminosities from parallaxes
- A symmetrical, well behaved error in parallax is transformed into an asymmetrical error in distance

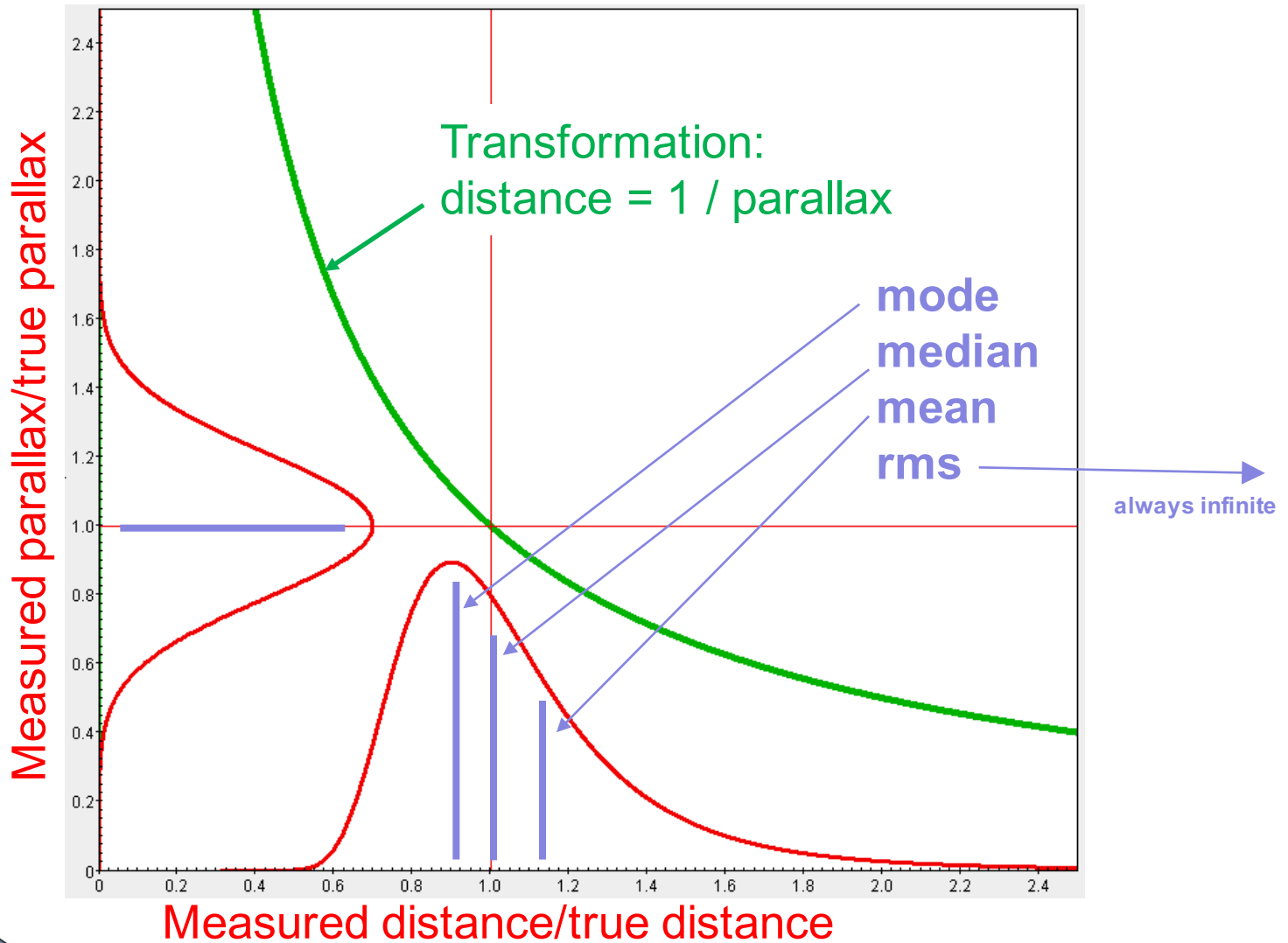
Error distribution comparison: star at 100pc and parallax error 2mas parallax and distance (non normalised)



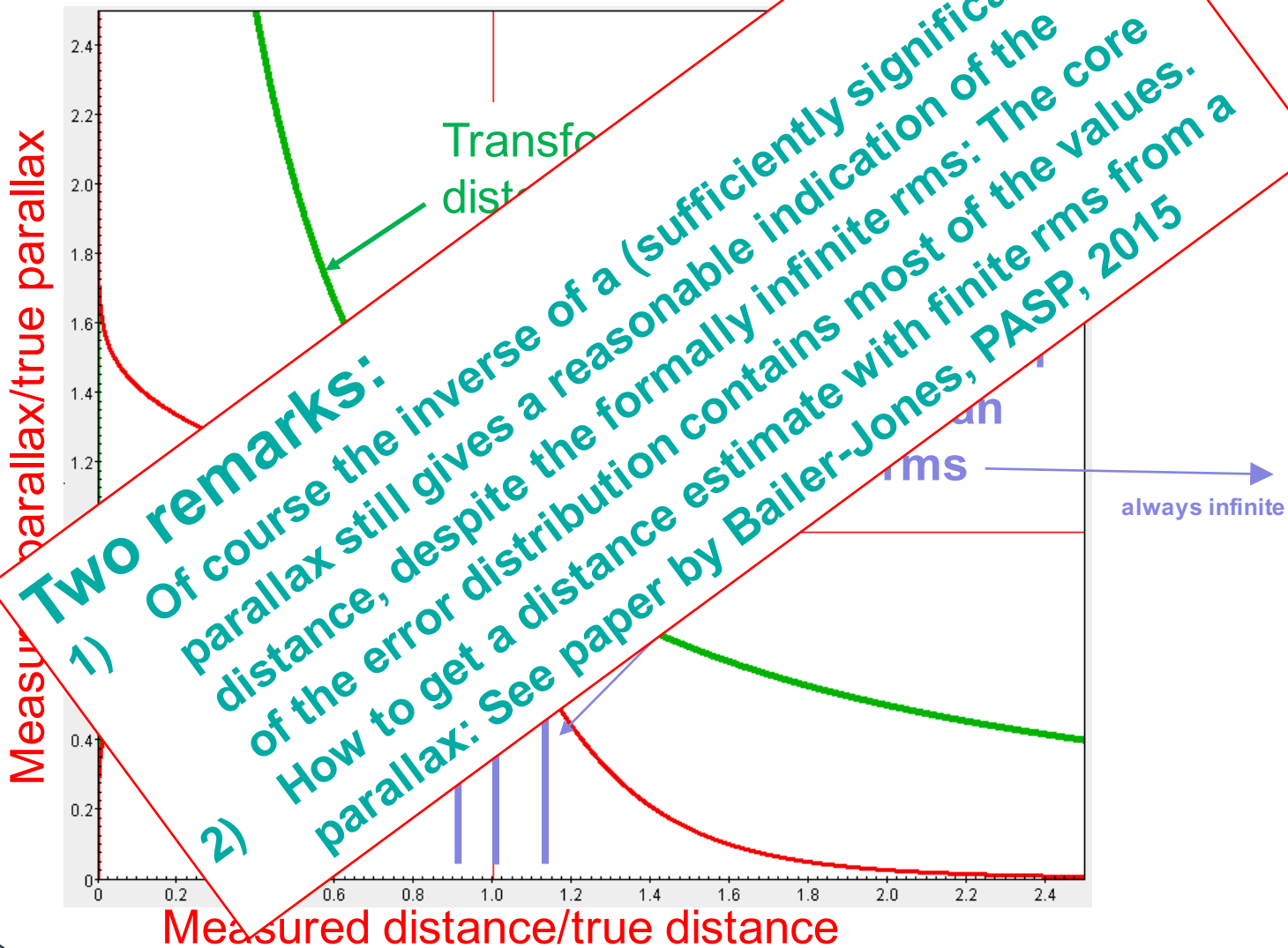
Error distribution comparison: parallax versus distance



Error distribution comparison: parallax versus distance



Error distribution comparison: parallax versus distance



Two remarks:

- 1) Of course the inverse of a (sufficiently significant) parallax still gives a reasonable indication of the distance, despite the formally infinite rms: The core of the error distribution contains most of the values.
- 2) How to get a distance estimate with finite rms from a parallax: See paper by Bailer-Jones, PASP, 2015



gaia

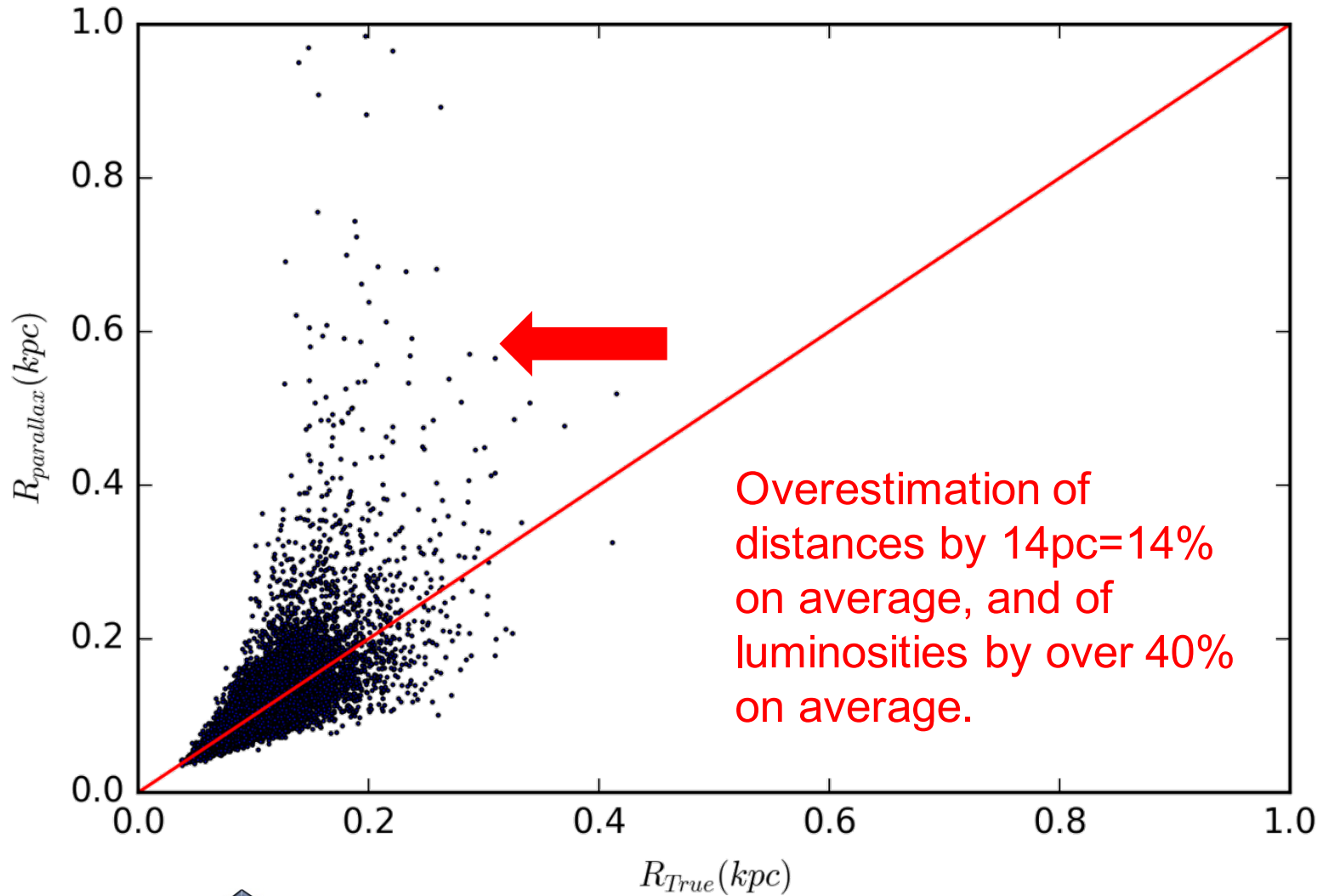


Gaia
DPAC
Data Processing & Analysis Consortium



Sample simulation with a parallax error of 2mas

True distance vs. distance from parallax



How to take this into account

- Avoid using transformations as much as possible
- If unavoidable:
 - Do fits in the plane of parallaxes (e.g. PL relations using ABL method*) where errors are well behaved
 - Do any averaging in parallaxes and then do the transformation (e.g. distance to an open cluster)
 - Always estimate the remaining effect (analytically or with simulations)

*Astrometry-Based Luminosity (ABL) method

$$a_V = 10^{0.2M_V} = \pi 10^{\frac{m_V + 5}{5}}$$

This quantity is:

- related to luminosity
(sqrt of inverse luminosity)
- a linear function of parallax
- thus nicely behaved
- thus can be averaged safely

Also beware of additional assumptions

- For instance about the absorption when calculating absolute magnitudes from parallaxes

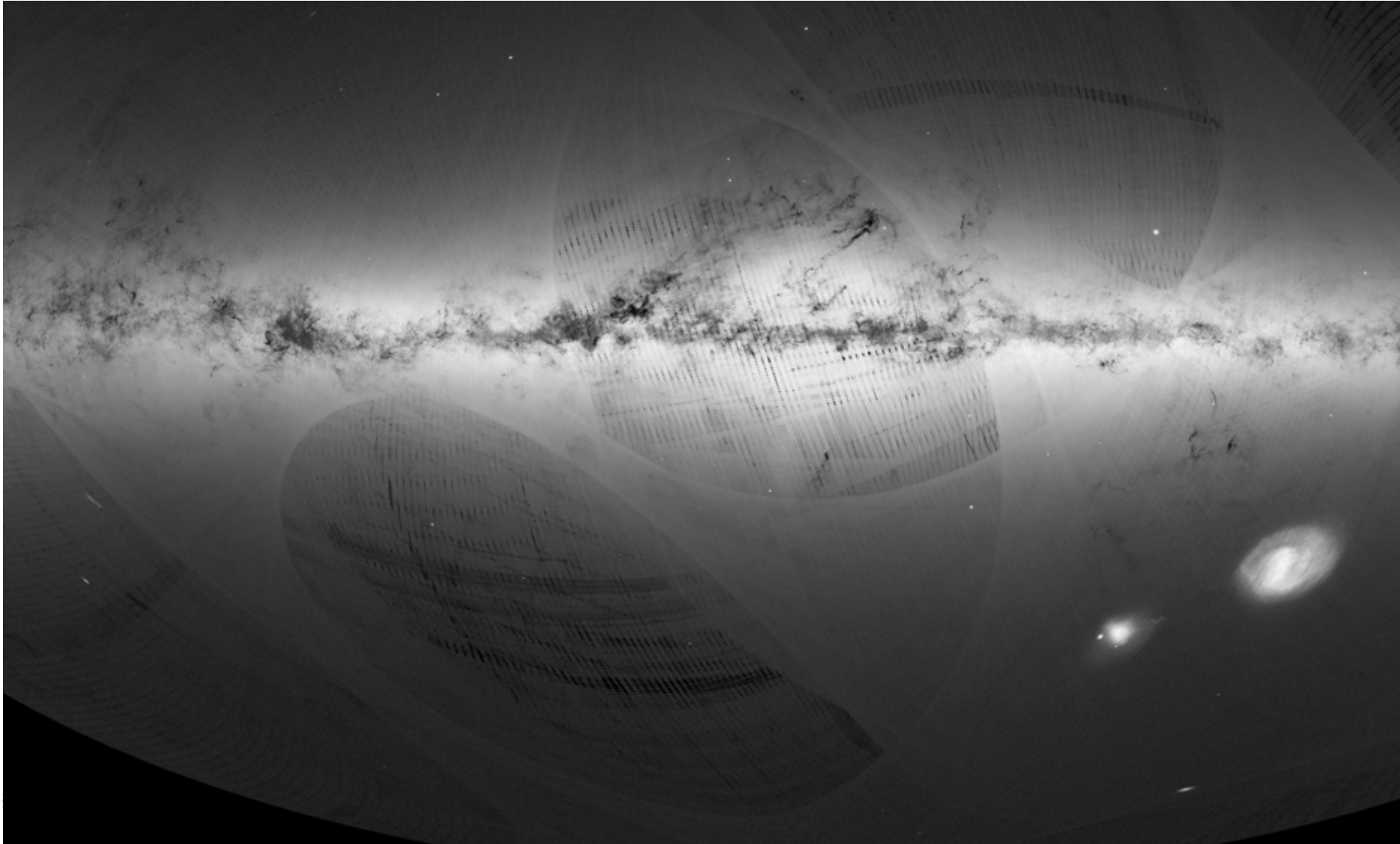
Chapter 5: Sample censorships

Completeness/representativeness:

we want to have the complete population of objects,
or at least a subsample which is representative for a given purpose

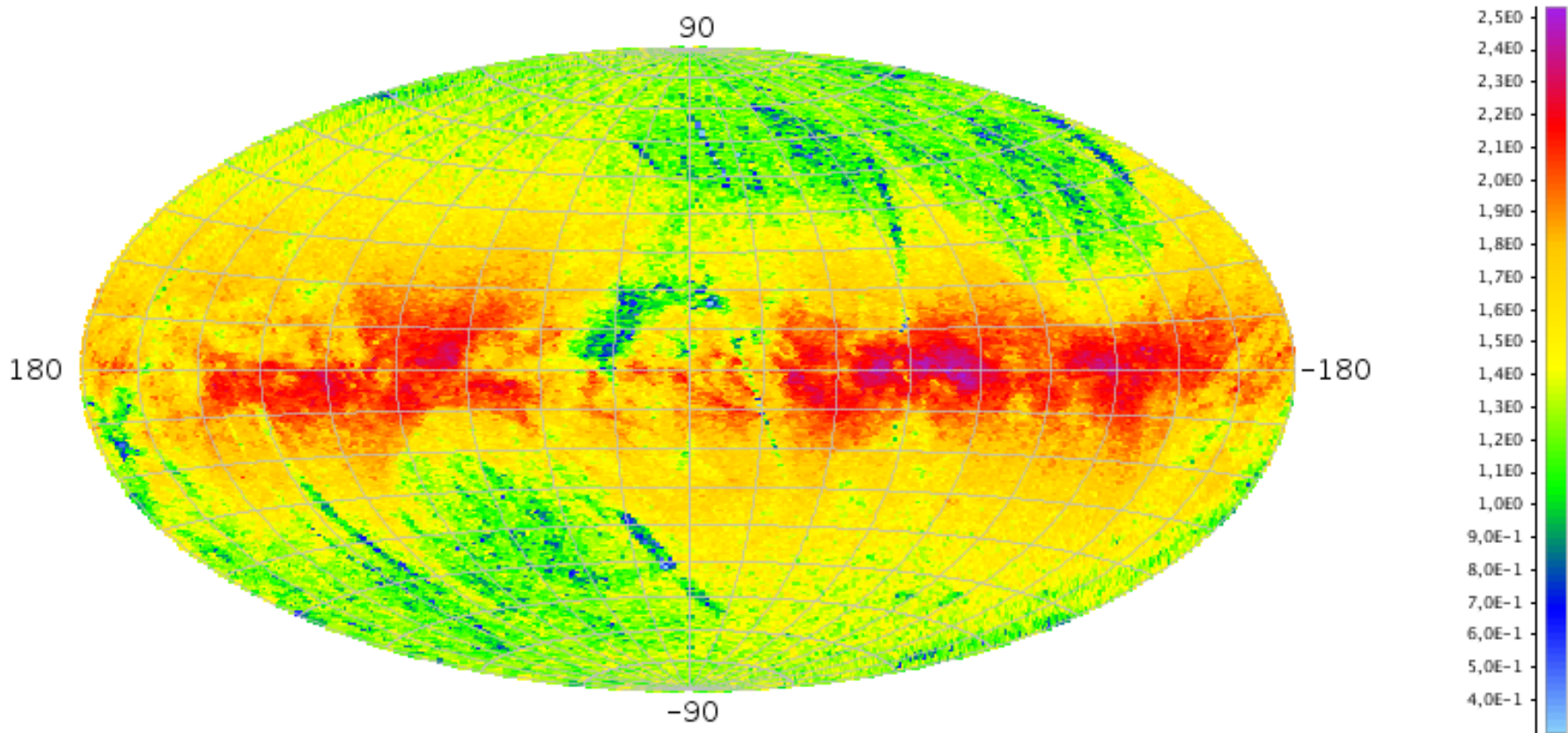
- This usually is not the case
- DR1 is a very complex dataset, its completeness or representativeness can not be guaranteed for any specific purpose

Example: Gaia DR1,
significant completeness variations
as a function of the sky position



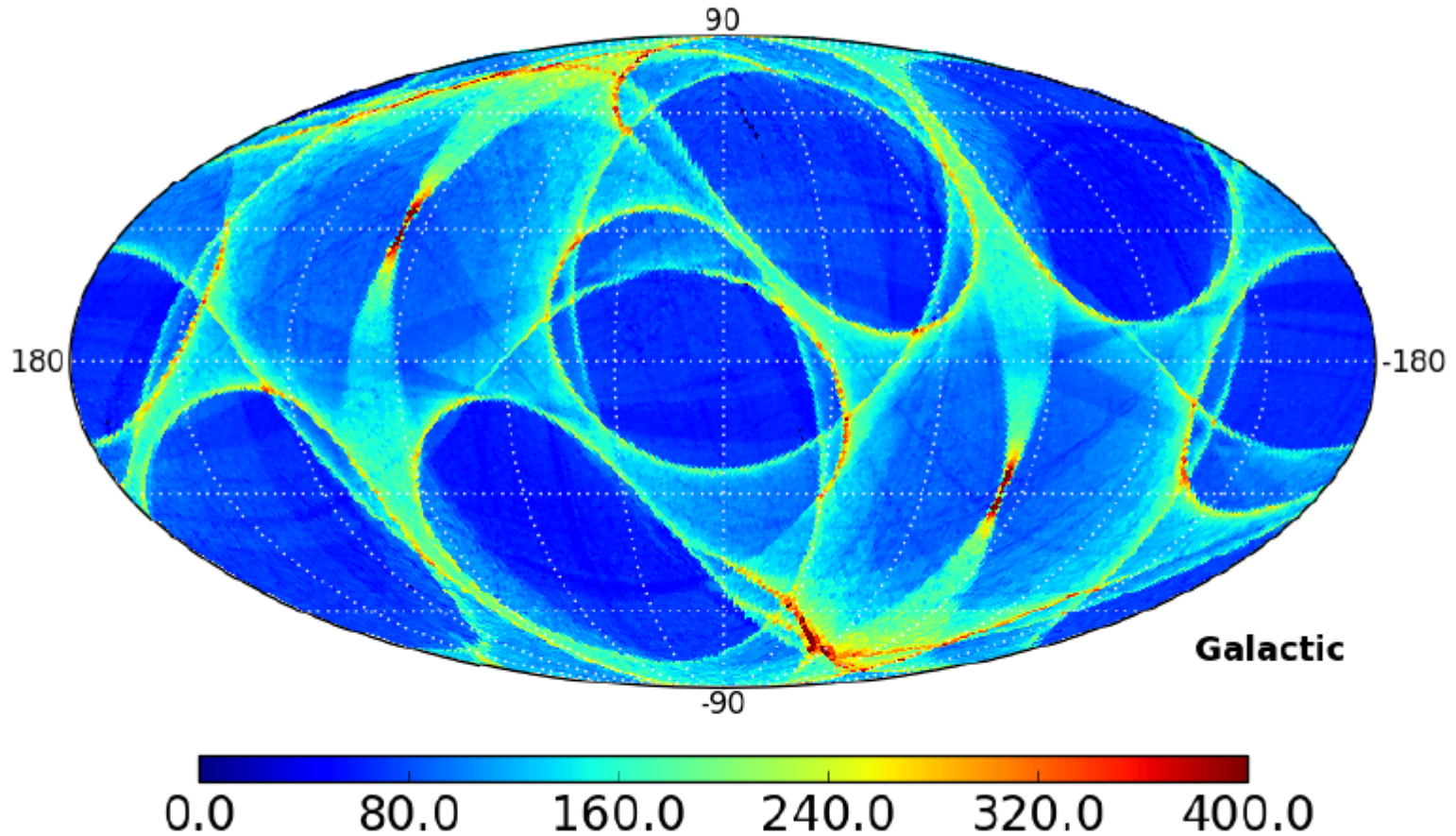
Significant completeness variations as a function of the sky position

Total log sky density in GAL coordinates (Log. of the number of objects). Objects: 2057050. Objects Out: 0

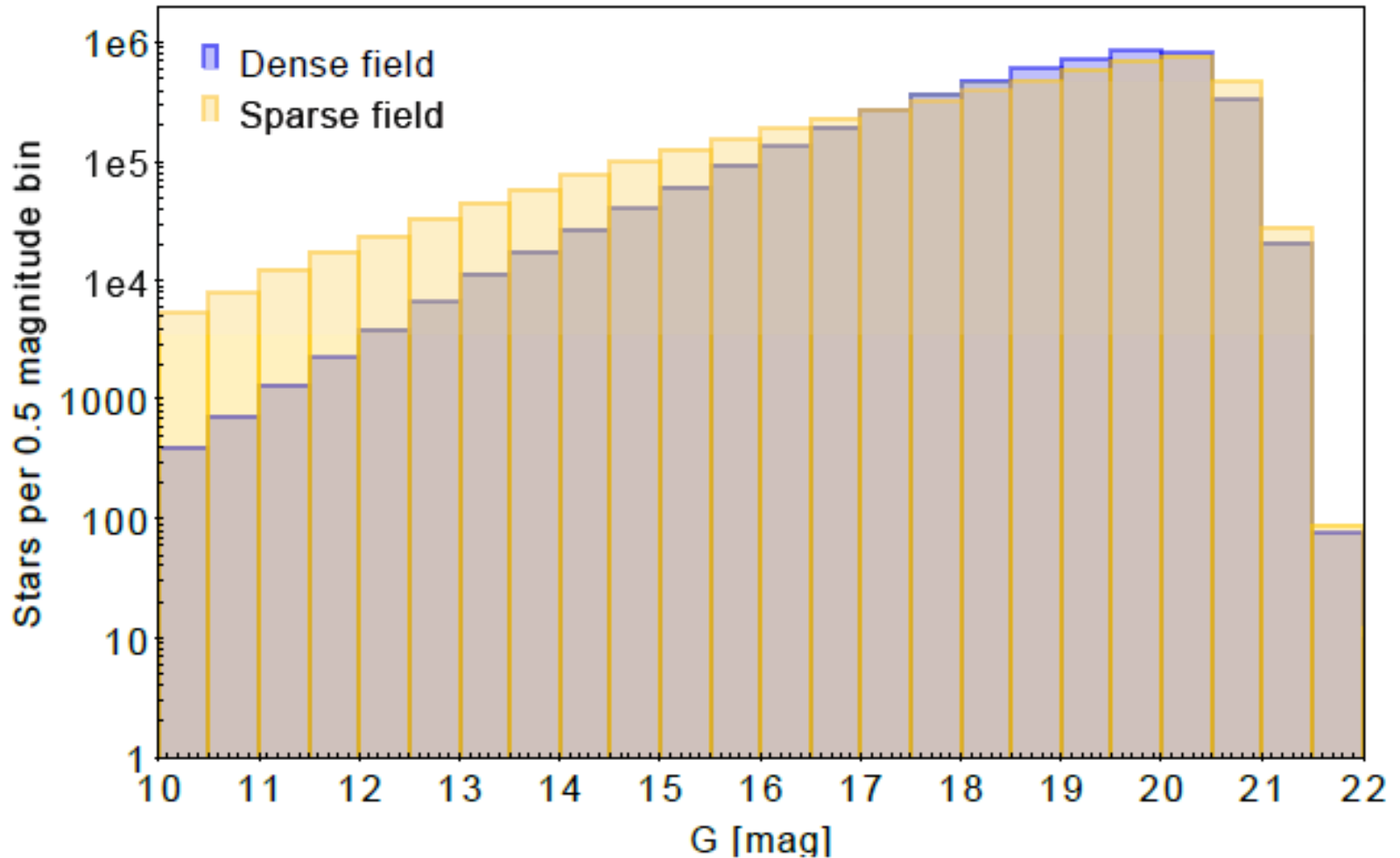


Complex selection of astrometry (e.g. Nobs)

TGAS Number of Good Observations Along Scan



Not complete in magnitude or color



How to take this into account

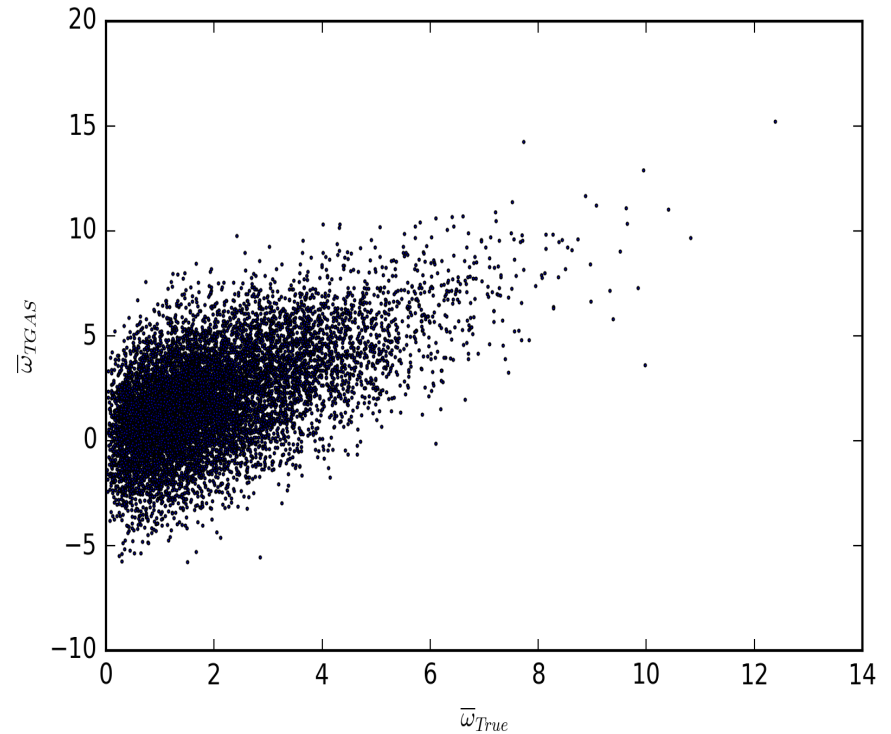
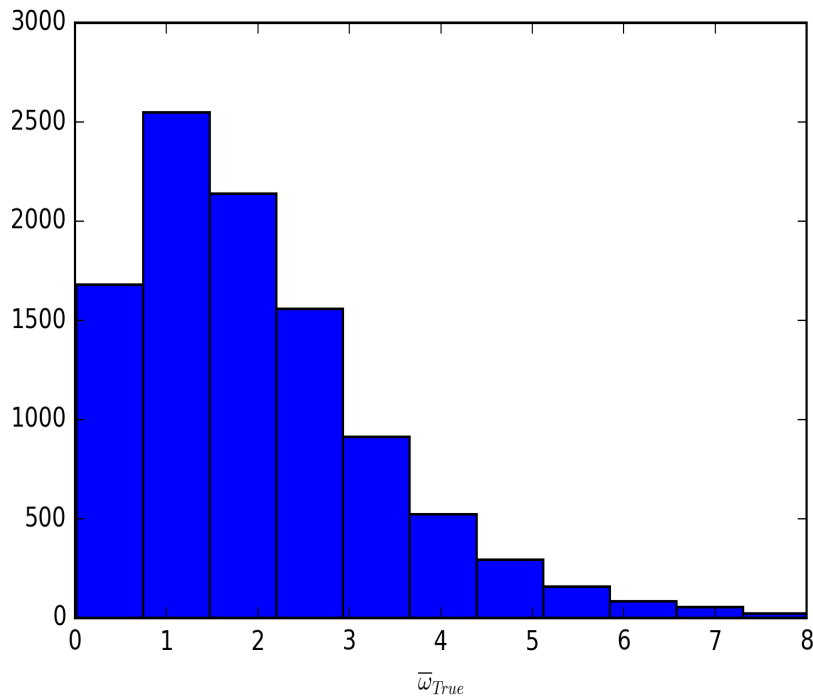
- **Very difficult**, will depend on your specific purpose
- Analyze if the problem exists, and try to determine if the known censorships are correlated with the parameter you are analyzing (see Gaia DR1 validation paper, A&A)
- If not possible analytically:
 - At least do some simulations to evaluate the possible effects

IMPORTANT: do not make things worse by adding your own additional censorships

- **This is especially important for parallaxes**
- Avoid removing negative parallaxes; this removes valid information, and it biases the sample for distant stars
- Avoid selecting subsamples on parallax relative error. This also removes information, and again it biases the sample for distant stars
- **Use instead fitting methods able to use all available data (e.g. Bayesian methods) and always work on the observables space (e.g. on parallaxes, not on distances or luminosities)**

Example: Original (complete) dataset

(assuming Gaussian errors in parallax of 2mas, and some “typical” true-distance distribution)



Simulation !

Average diff. of parallaxes = 0.002 mas

fine and expected: within $2\text{mas}/\sqrt{10000}$

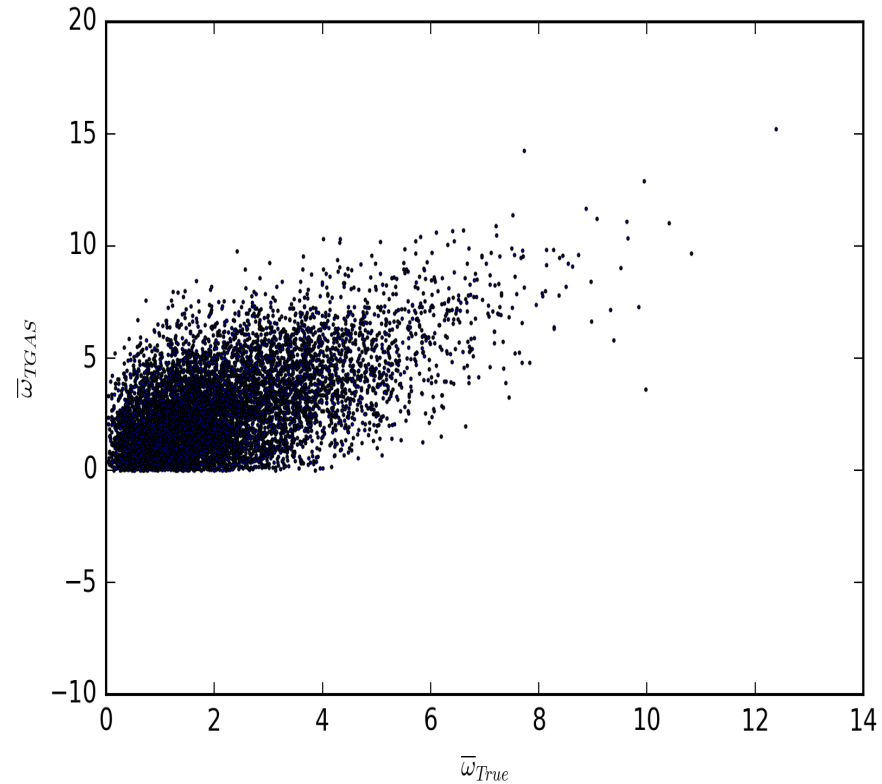
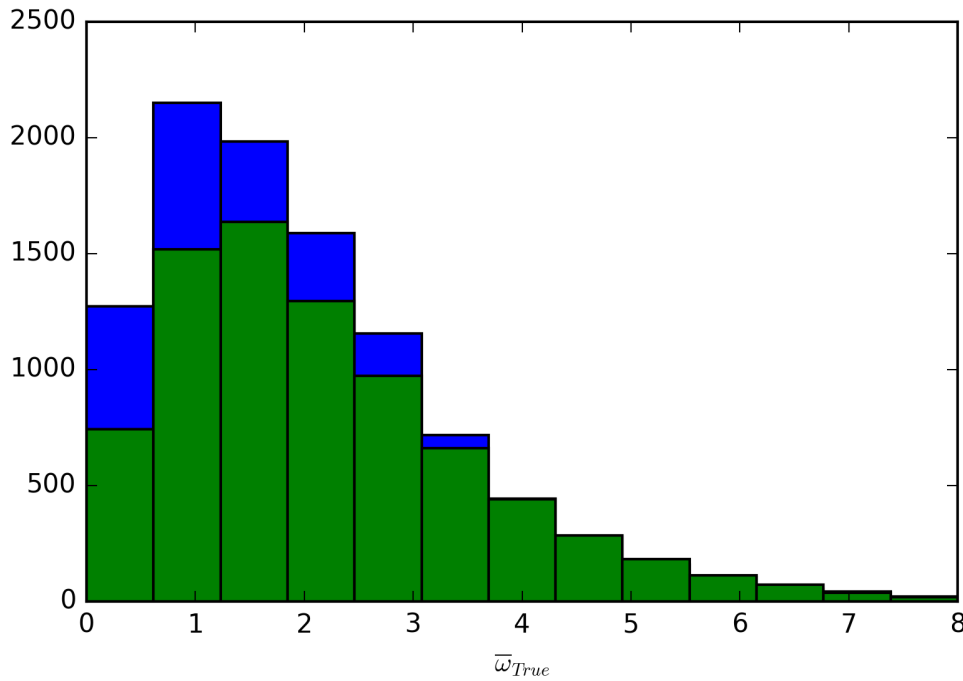


gaia



Example: removing negative parallaxes

Favours large parallaxes



Simulation!

Average diff. of parallaxes = 0.65 mas

disastrous

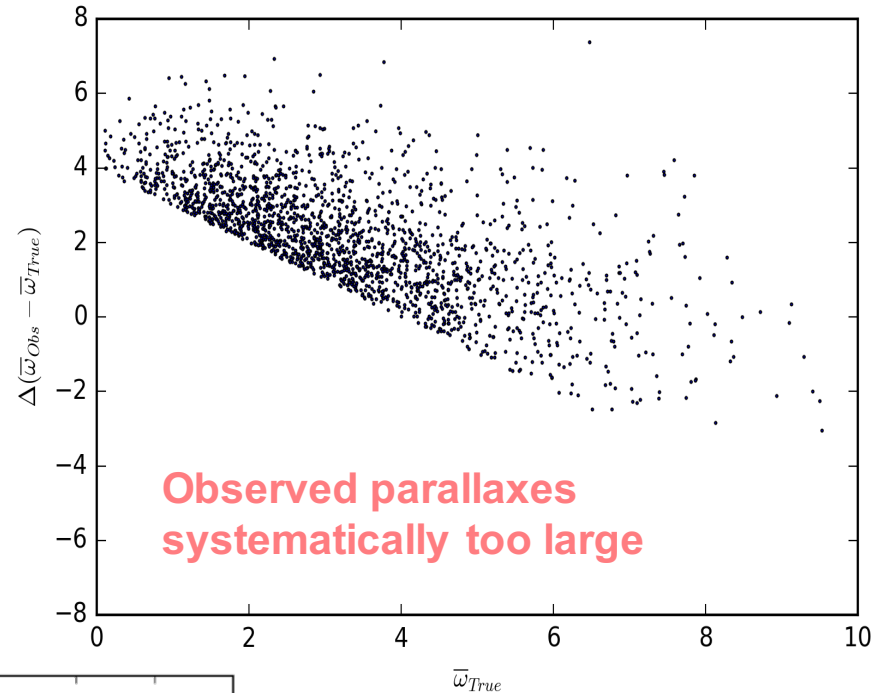
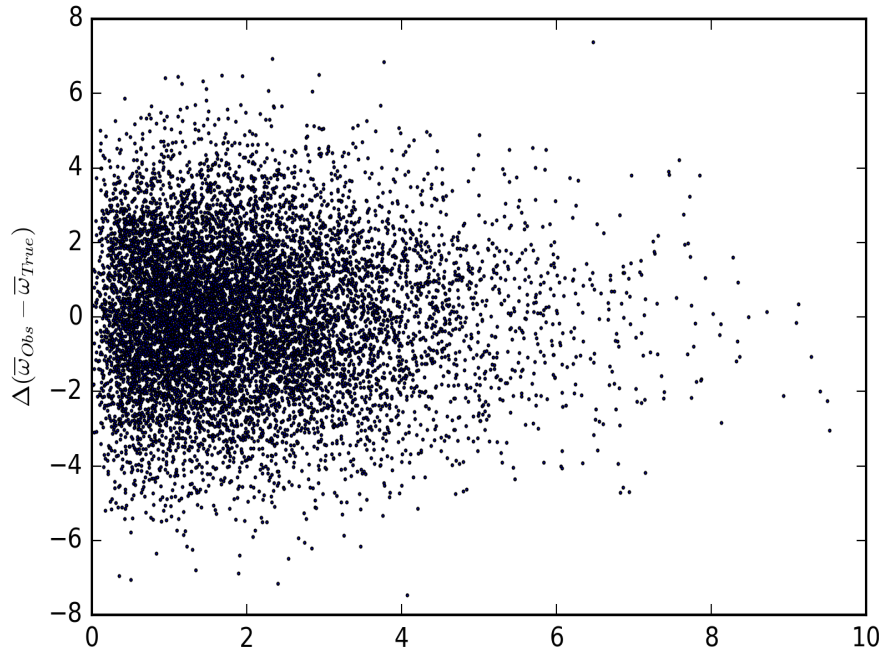


gaia

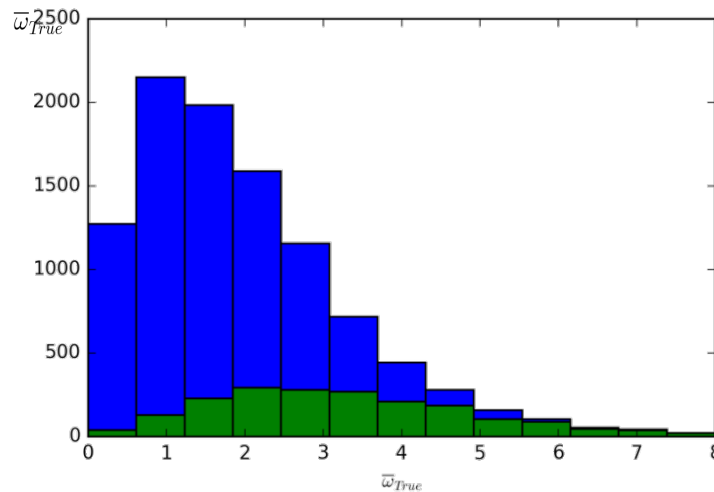


Example: removing $\sigma_{\text{Par}}/\text{Par} > 50\%$

Favours errors making parallax larger



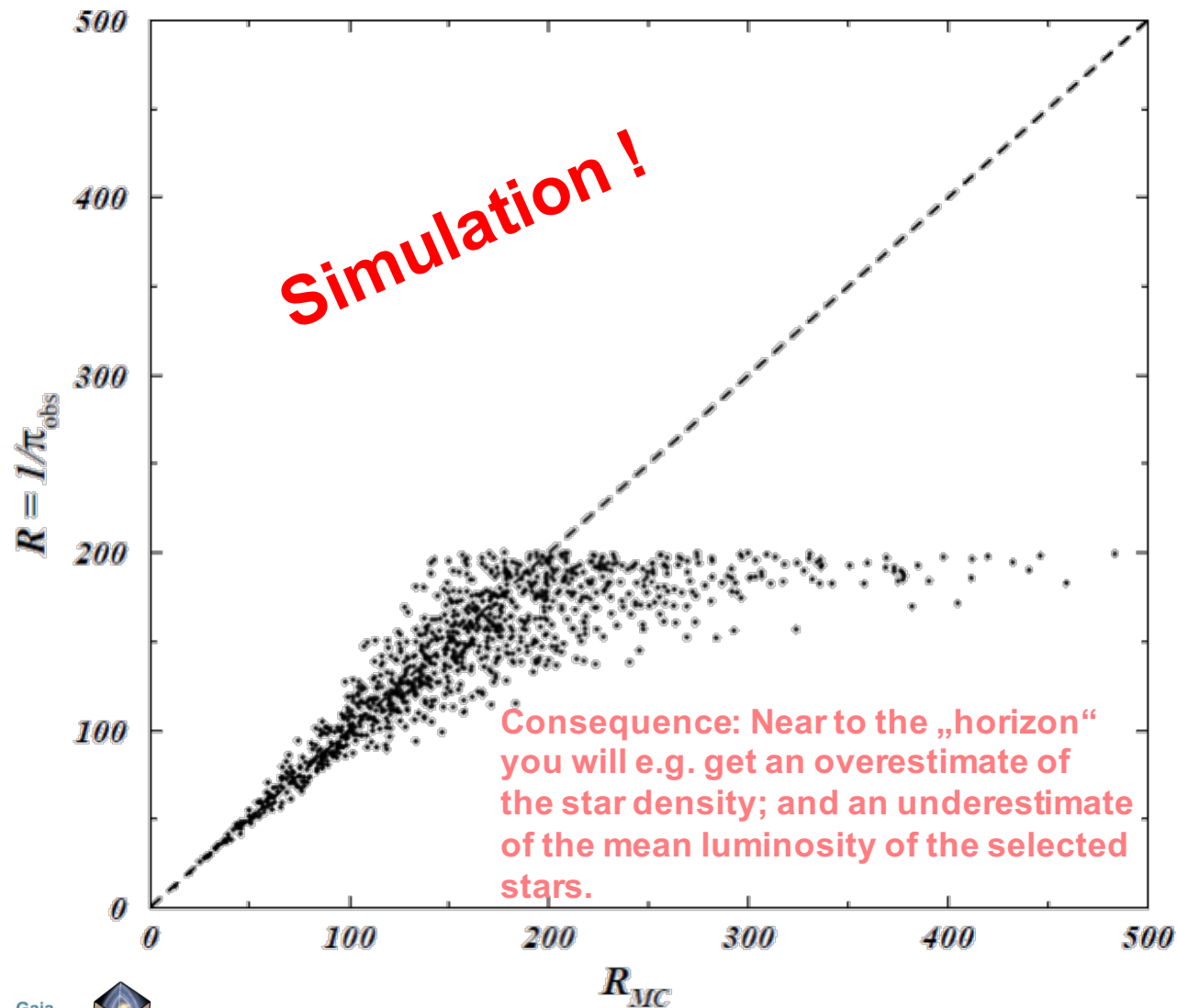
Simulation!



Average diff. of
parallaxes = 2.2 mas

Example: truncation by observed parallax

Favours objects at large distances (small true parallax)



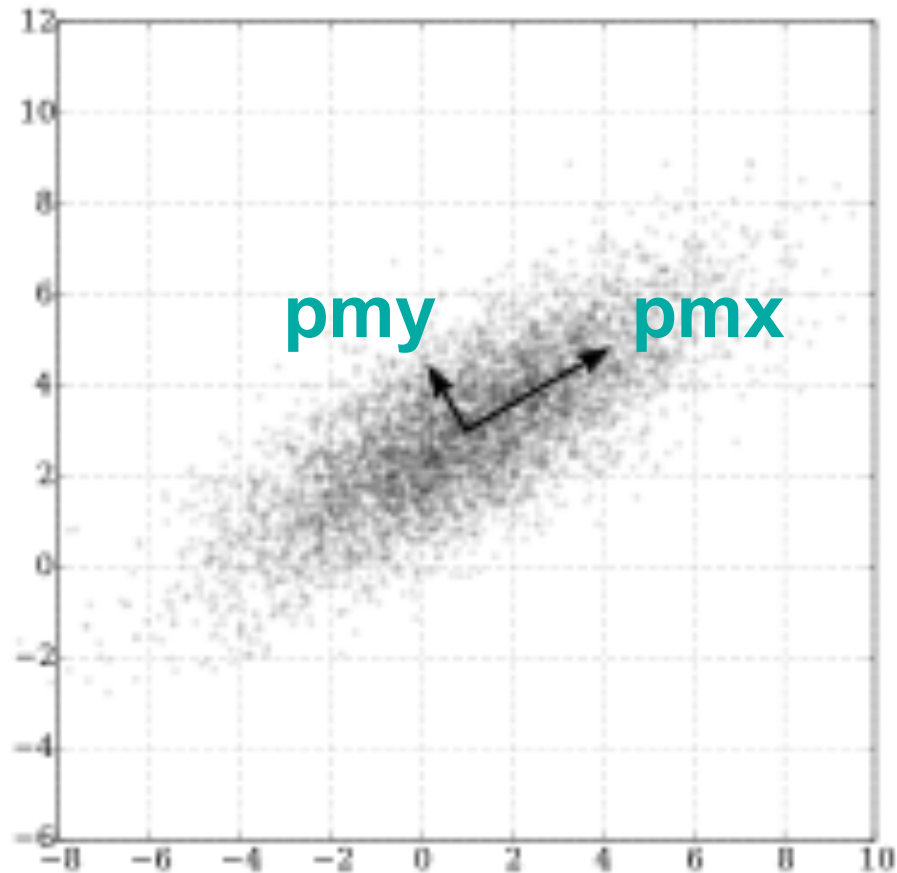
gaia



Thank you

Appendix

Uncorrelated quantities from correlated catalogue values



Given:

pma , pmd ,
 $\sigma(pma)$, $\sigma(pmd)$, $\text{corr}(pma, pmd)$

Wanted: orientation and principal axes of the error ellipse

Go to rotated coordinate system x, y . The two proper-motion components pmx and pmy are uncorrelated:

$$pmx = pmd \cdot \cos(\theta) + pma \cdot \sin(\theta)$$

$$pmy = -pmd \cdot \sin(\theta) + pma \cdot \cos(\theta)$$

Question:

Which θ ?

And which $\sigma(pm_x)$, $\sigma(pm_y)$?

Uncorrelated quantities from correlated catalogue values

Keyword: Eigenvalue decomposition (of the relevant covariance matrix part)

Example for the “looks” of a covariance matrix (2×2 , proper motions only):

$$\begin{pmatrix} \sigma_{\mu_{\alpha^*}}^2 & \text{COV}_{\mu_{\alpha^*}, \mu_{\delta}} \\ \text{COV}_{\mu_{\alpha^*}, \mu_{\delta}} & \sigma_{\mu_{\delta}}^2 \end{pmatrix}$$

Note: $\text{cov}_{\mu_{\alpha^*}, \mu_{\delta}} = \text{corr}_{\mu_{\alpha^*}, \mu_{\delta}} \sigma_{\mu_{\alpha^*}} \sigma_{\mu_{\delta}}$

Solution of the Eigenvalue decomposition for 2 dimensions:

The maxima and minima of the variance (the eigenvalues of the matrix) are:

$$\begin{aligned} \sigma_{\mu_x}^2 &= \frac{1}{2} \left(\sigma_{\mu_{\alpha^*}}^2 + \sigma_{\mu_{\delta}}^2 + \sqrt{[\sigma_{\mu_{\alpha^*}}^2 + \sigma_{\mu_{\delta}}^2]^2 - 4\text{cov}_{\mu_{\alpha^*}, \mu_{\delta}}^2} \right) \\ \sigma_{\mu_y}^2 &= \frac{1}{2} \left(\sigma_{\mu_{\alpha^*}}^2 + \sigma_{\mu_{\delta}}^2 - \sqrt{[\sigma_{\mu_{\alpha^*}}^2 + \sigma_{\mu_{\delta}}^2]^2 - 4\text{cov}_{\mu_{\alpha^*}, \mu_{\delta}}^2} \right) \\ \tan(\theta) &= \frac{\sigma_{\mu_{\alpha^*}}^2 - \sigma_{\mu_{\delta}}^2}{\text{COV}_{\mu_{\alpha^*}, \mu_{\delta}}} \end{aligned}$$

Note 1: the $\pm 180^\circ$ ambiguity of the tangent does not matter in this case

Note 2: for $\text{cov}_{\mu_{\alpha^*}, \mu_{\delta}} = 0$, then $\theta = 0$ if $\sigma_{\mu_{\delta}} > \sigma_{\mu_{\alpha^*}}$, else $\theta = 90^\circ$ and the values are trivial

Even more tedious formulae for 3 dimensions; better use matrix routines for 3d and higher dimensions.

Uncorrelated quantities from correlated catalogue values

Keyword: Eigenvalue decomposition
(of the relevant covariance matrix part)

Example for the “looks” of a covariance matrix (2 by 2, proper motions only):

$$\left\{ \begin{array}{cc} \sigma^2(\text{pma}) & \text{cov}(\text{pma}, \text{pmd}) \\ \text{cov}(\text{pma}, \text{pmd}) & \sigma^2(\text{pmd}) \end{array} \right\}$$

Note: $\text{cov}(\text{pma}, \text{pmd}) = \text{corr}(\text{pma}, \text{pmd}) * \sigma(\text{pma}) * \sigma(\text{pmd})$

Solution of the Eigenvalue decomposition for 2 dimensions: (promised during the talk to be added here)

The maxima and minima of the variance (the eigenvalues of the matrix) are:

$$\sigma^2(\text{pmx}) = 1/2 * (\sigma^2(\text{pma}) + \sigma^2(\text{pmd}) + \sqrt{(\sigma^2(\text{pma}) + \sigma^2(\text{pmd}))^2 - 4\text{cov}^2(\text{pma}, \text{pmd})})$$

$$\sigma^2(\text{pmy}) = 1/2 * (\sigma^2(\text{pma}) + \sigma^2(\text{pmd}) - \sqrt{(\sigma^2(\text{pma}) + \sigma^2(\text{pmd}))^2 - 4\text{cov}^2(\text{pma}, \text{pmd})})$$

$\tan(\theta) = (\sigma^2(\text{pma}) - \sigma^2(\text{pmd})) / \text{cov}(\text{pma}, \text{pmd})$; note 1: the +/- 180 deg ambiguity of the tangens does not matter in this case.
note 2: for $\text{cov}(\text{pma}, \text{pmd})=0$, then $\theta=0$ if $\sigma(\text{pmd}) > \sigma(\text{pma})$, else $\theta=90\text{deg}$, and the values are trivial

Even more tedious formulae for 3 dimensions; better use matrix routines for 3d and higher dimensions.

Sorry for the clumsy formula notation, but I didn't find the time to typeset them more nicely. Volunteers are invited to email me ☺

During this presentation

- about 1 million stars were measured by Gaia,
- roughly 10 million astrometric measurements were taken,
- about 300,000 spectra were made of 100,000 stars