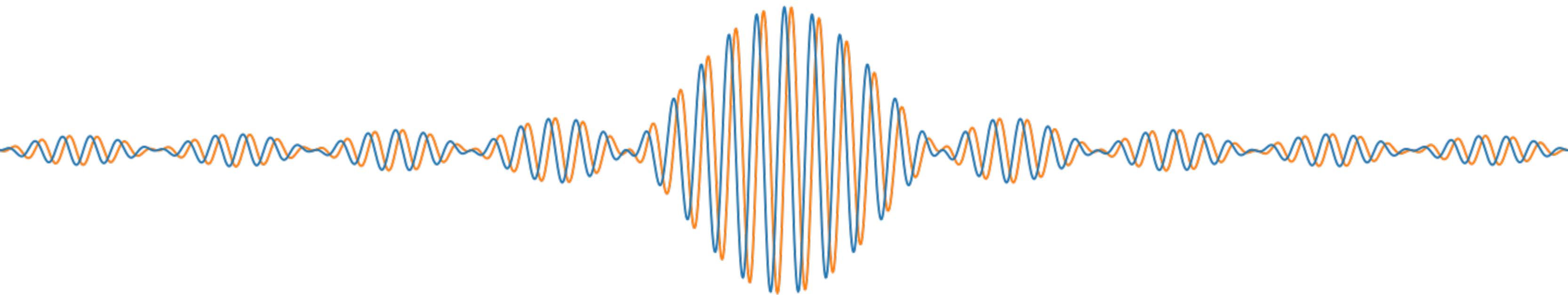


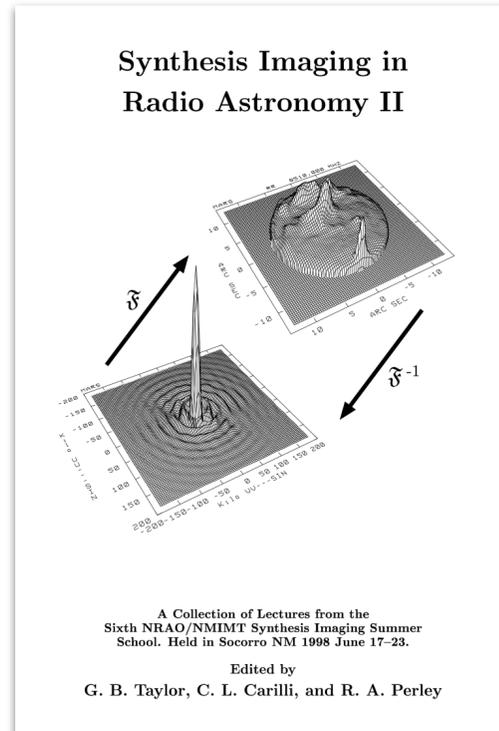
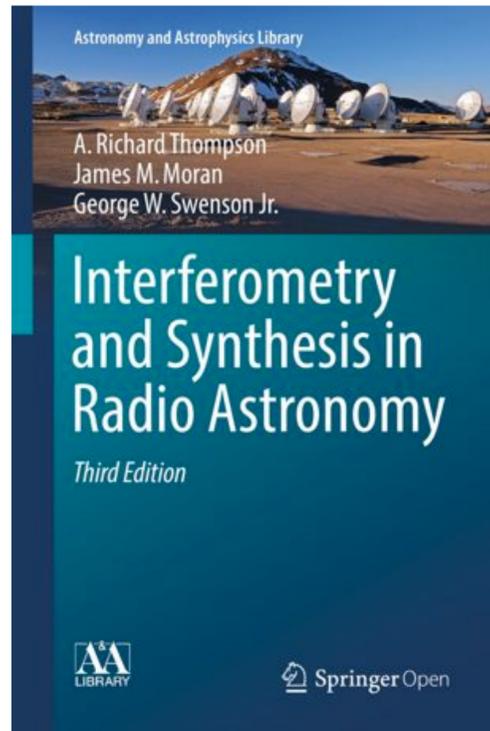


VLBI Data & Errors



Lindy Blackburn

*Event Horizon Telescope Collaboration
Center for Astrophysics | Harvard & Smithsonian*



Tetsuo Sasao and André B. Fletcher
Introduction to VLBI Systems
Chapter 4
Lecture Notes for KVN Students
Partly based on Ajou University Lecture Notes
(to be further edited)
Version 1. (Unfinished)
Issued on February 19, 2006.

Very Long Baseline Interferometry

Contents

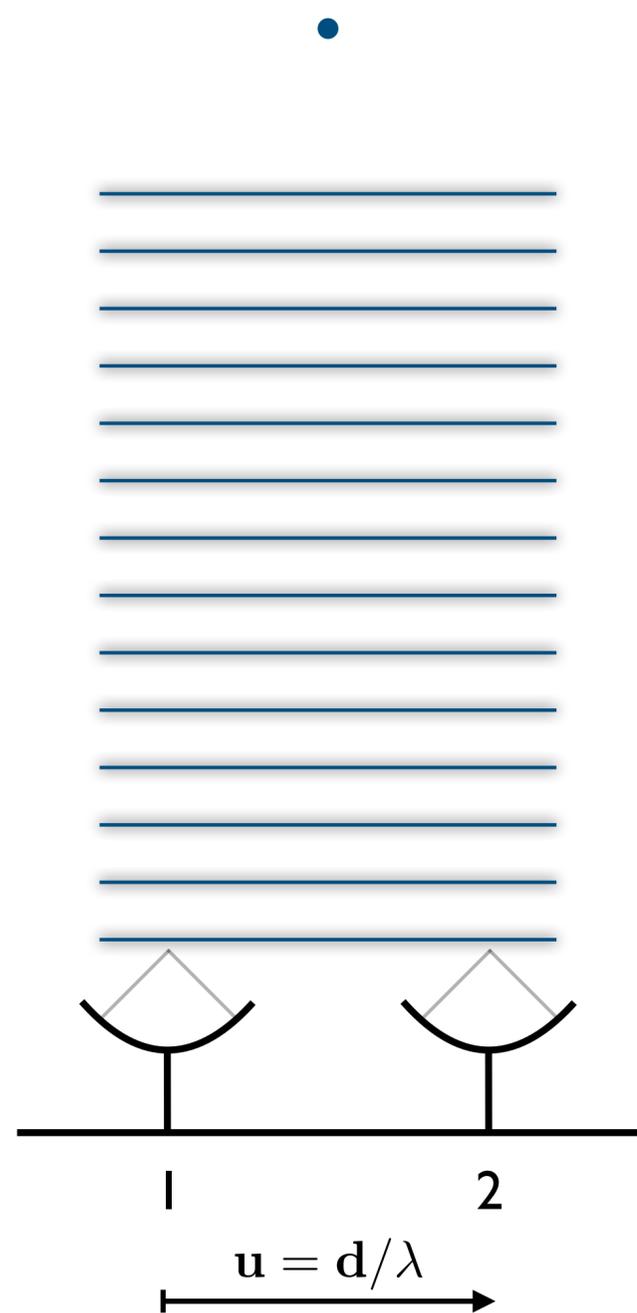
1 Technologies Which Made VLBI Possible	4
1.1 Basics of Digital Data Processing	5
1.1.1 Analog Processing Versus Digital Processing	5
1.1.2 Sampling and Clipping	6
1.1.3 Discrete-Time Random Process	6
1.1.4 Stationary Discrete-Time Random Process	9
1.1.5 Sampling	10
1.1.6 Comb Function	12
1.1.7 Fourier Transform of a Comb Function Is a Comb Function	13
1.1.8 Spectra of Discrete-Time Processes	14
1.1.9 Spectra of Sampled Data	17
1.1.10 Inverse Relations for Spectra of Sampled Data	19
1.1.11 Sampling Theorem	20
1.1.12 Optimum Sampling Interval	22
1.1.13 Sampling Function	25
1.1.14 Correlations of Nyquist Sampled Data with Rectangular Passband Spectra	28
1.1.15 S/N Ratio of Correlator Output of Sampled Data	33
1.1.16 Nyquist Theorem and Nyquist Interval	36
1.1.17 Higher-Order Sampling in VLBI Receiving Systems	37
1.1.18 Clipping (or Quantization) of Analog Data	39
1.1.19 Probability Distribution of Clipped Data	42
1.1.20 Cross-Correlation of 1-Bit Quantized Data:	
van Vleck Relationship	46
1.1.21 van Vleck Relationship in Autocorrelation	51

Lecture	Speaker
Basics of Radio Astronomy	Like Young (NRT)
Antennas & Receivers in Radio Astronomy	Mark McKinnon (NRAO)
Fundamentals of Radio Interferometry I	Rick Perley (NRAO)
Fundamentals of Radio Interferometry II	Rick Perley (NRAO)
Fundamentals of Radio Interferometry III	Rick Perley (NRAO)
Interferometry of Solar System Objects	Bryan Butler (NRAO)
Cross Correlators	Adam Deller (ASTRON)
The High-Redshift Universe, Magnified	Dan Marrone (UA)
Calibration	George Moellenbrock (NRAO)
Imaging and Deconvolution	David Wilner (CRA)
Advanced Calibration I	Crystal Bragan (NRAO)
Advanced Calibration II	Crystal Bragan (NRAO)
The XLA Sky Survey	Claire Chandler (NRAO)
Multi-messenger Exploration of the Transient Sky	Alexandra Conn (TTU)
Spectral Line Data Analysis	Vivus Pridmore (LJMU)
Polarization	Frank Schinzel (NRAO)
Masking	Brian Mason (NRAO)
Very Long Baseline Interferometry	Adam Deller (ASTRON)
Low-Frequency Interferometry	Tracy Clarke (MNL)
Astrochemistry	Brett McGuire (NRAO)
Wideband and Wide-field Imaging I	Unvashi Rao (NRAO)
Wideband and Wide-field Imaging II	Unvashi Rao (NRAO)
Protoplanetary Disks	Paul U. LAMN
II Zw 40: A Test Case for Studying Barium Cycling in the Nearby Universe	Aravinda Kigley (NRAO)
Error Recognition	Greg Taylor (LJMU)
Image and Non-Imaging Analysis	Greg Taylor (LJMU)
Array Design	Craig Walker (NRAO)
VLA Planning your Observation - Lecture	Lorant Spowerman (NRAO)
ALMA Planning your Observation - Lecture	Rachel Frisen (NRAO)

- Interferometry and Synthesis in Radio Astronomy
Thomson, Moran, Swenson (TMS)
- Synthesis Imaging in Radio Astronomy II
Ed: Taylor, Carilli, Perley
- KVN Lecture Notes
Sasao, Fletcher
- Synthesis and Imaging Workshop 2018 Presentations
NRAO

Van Cittert-Zernike theorem

- Relates spatial coherence of wavefront with brightness distribution of distant source

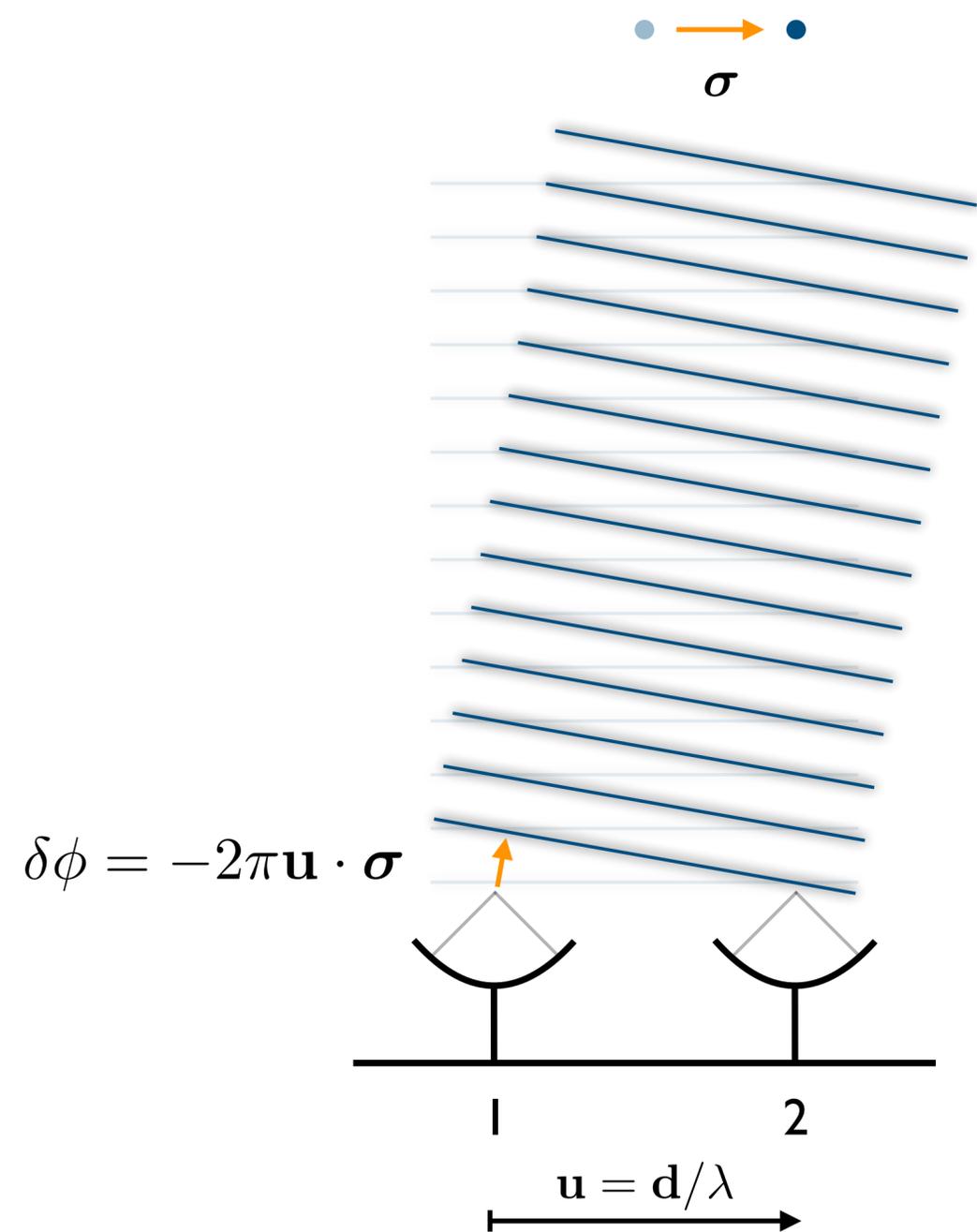


point source

$$\langle E_1 E_2^* \rangle = S_\nu$$

Van Cittert-Zernike theorem

- Relates spatial coherence of wavefront with brightness distribution of distant source



point source

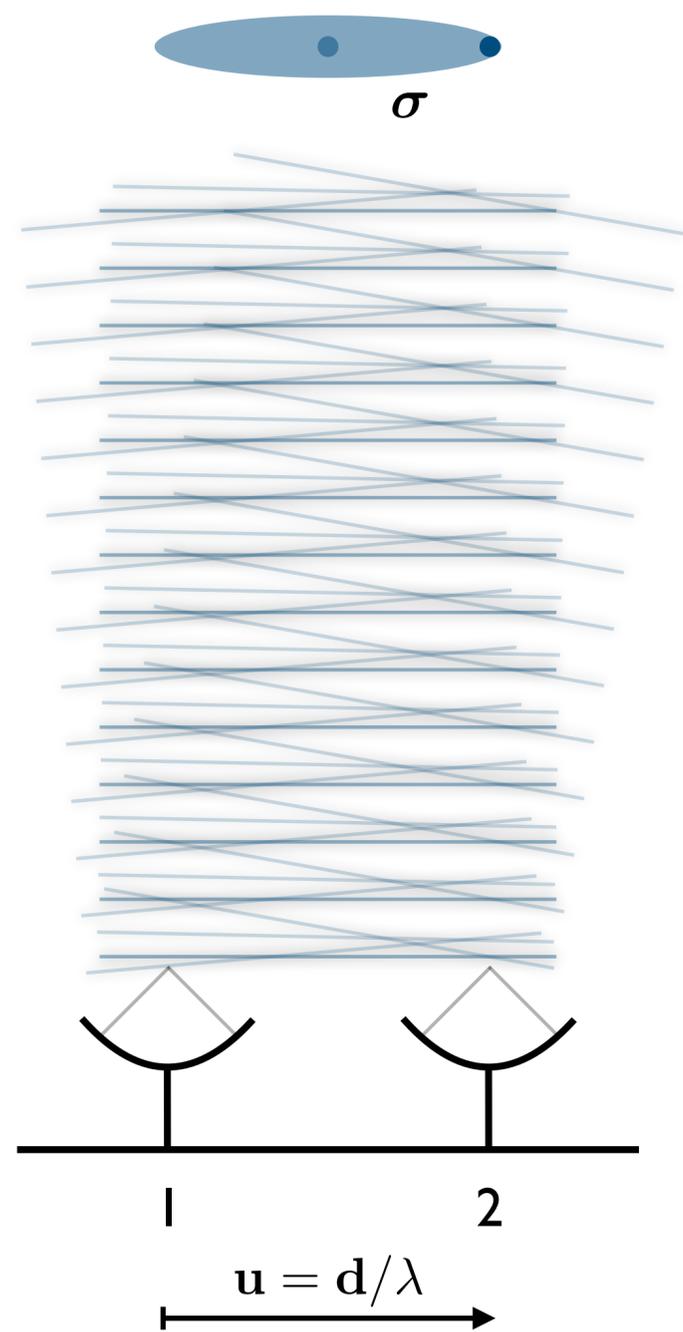
$$\langle E_1 E_2^* \rangle = S_\nu$$

shifted point source

$$\langle E_1 E_2^* \rangle = e^{-2\pi\mathbf{u} \cdot \boldsymbol{\sigma}} S_\nu$$

Van Cittert-Zernike theorem

- Relates spatial coherence of wavefront with brightness distribution of distant source



point source

$$\langle E_1 E_2^* \rangle = S_\nu$$

shifted point source

$$\langle E_1 E_2^* \rangle = e^{-2\pi \mathbf{u} \cdot \boldsymbol{\sigma}} S_\nu$$

extended source (integration over many point sources)

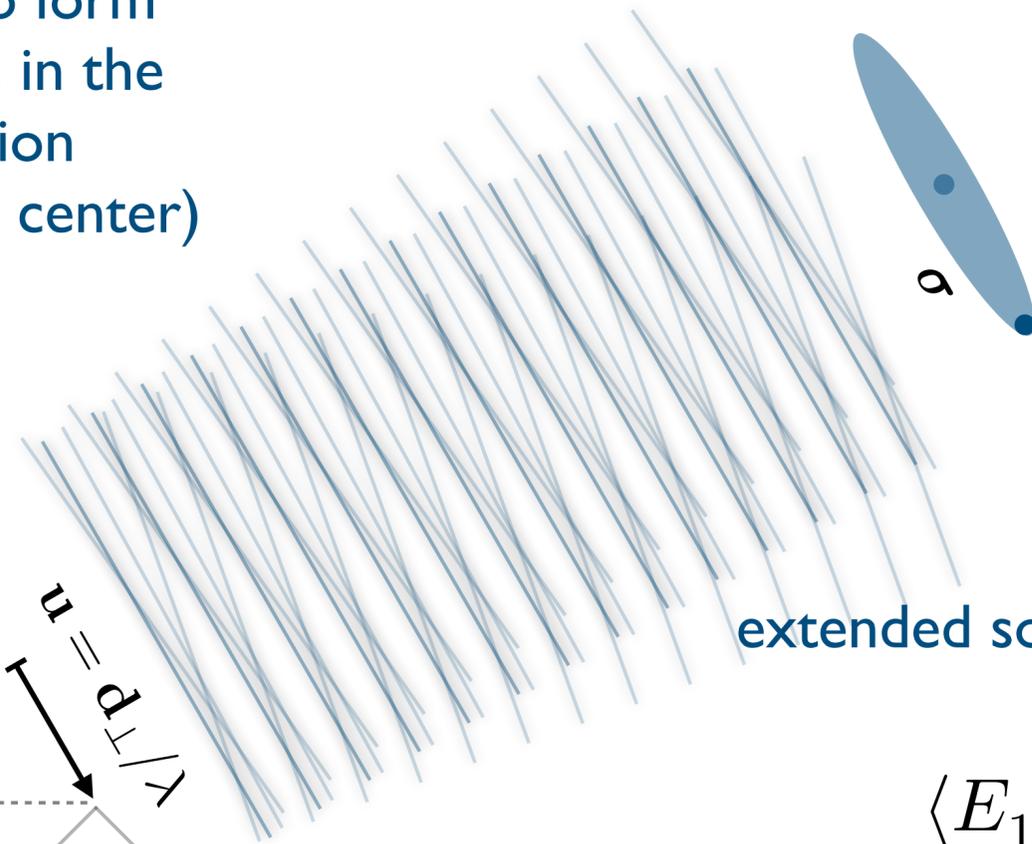
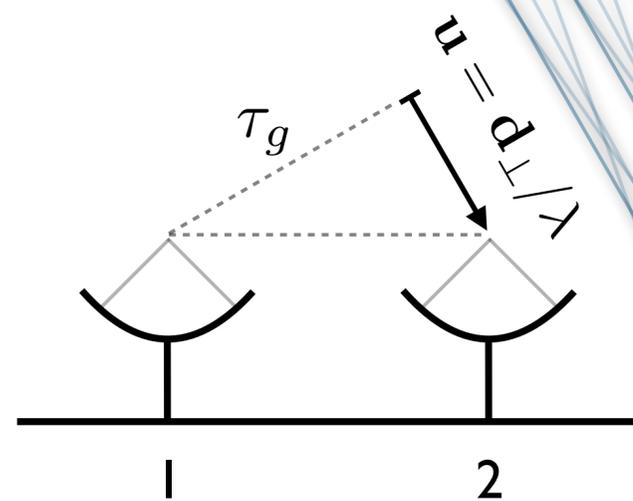
$$\begin{aligned} \langle E_1 E_2^* \rangle &= \iint e^{-2\pi \mathbf{u} \cdot \boldsymbol{\sigma}} I_\nu(\boldsymbol{\sigma}) d\Omega \\ &= \mathcal{V}(\mathbf{u}) \end{aligned}$$

“Visibility function”
encodes 2D complex spatial frequency
components of the sky brightness

Van Cittert-Zernike theorem

- Relates spatial coherence of wavefront with brightness distribution of distant source

Time-shift antenna to form baseline vector taken in the plane of propagation
(Linearize about phase center)



point source

$$\langle E_1 E_2^* \rangle = S_\nu$$

shifted point source

$$\langle E_1 E_2^* \rangle = e^{-2\pi \mathbf{u} \cdot \boldsymbol{\sigma}} S_\nu$$

extended source (integration over many point sources)

$$\begin{aligned} \langle E_1 E_2^* \rangle &= \iint e^{-2\pi \mathbf{u} \cdot \boldsymbol{\sigma}} I_\nu(\boldsymbol{\sigma}) d\Omega \\ &= \mathcal{V}(\mathbf{u}) \end{aligned}$$

“Visibility function”
encodes 2D complex spatial frequency
components of the sky brightness

Ingredients of a VLBI measurement



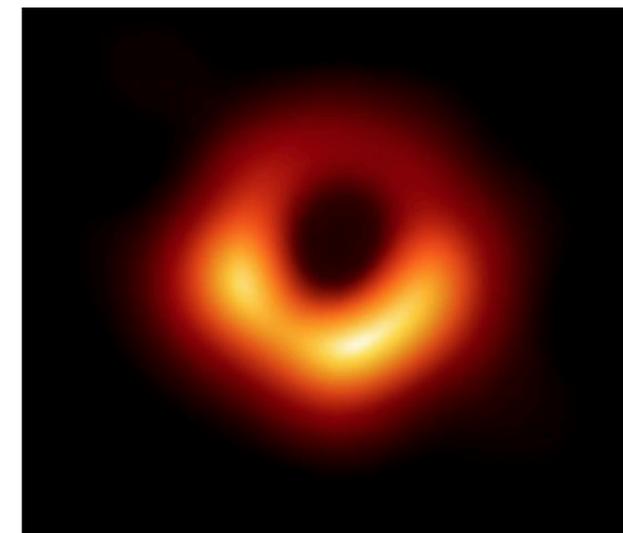
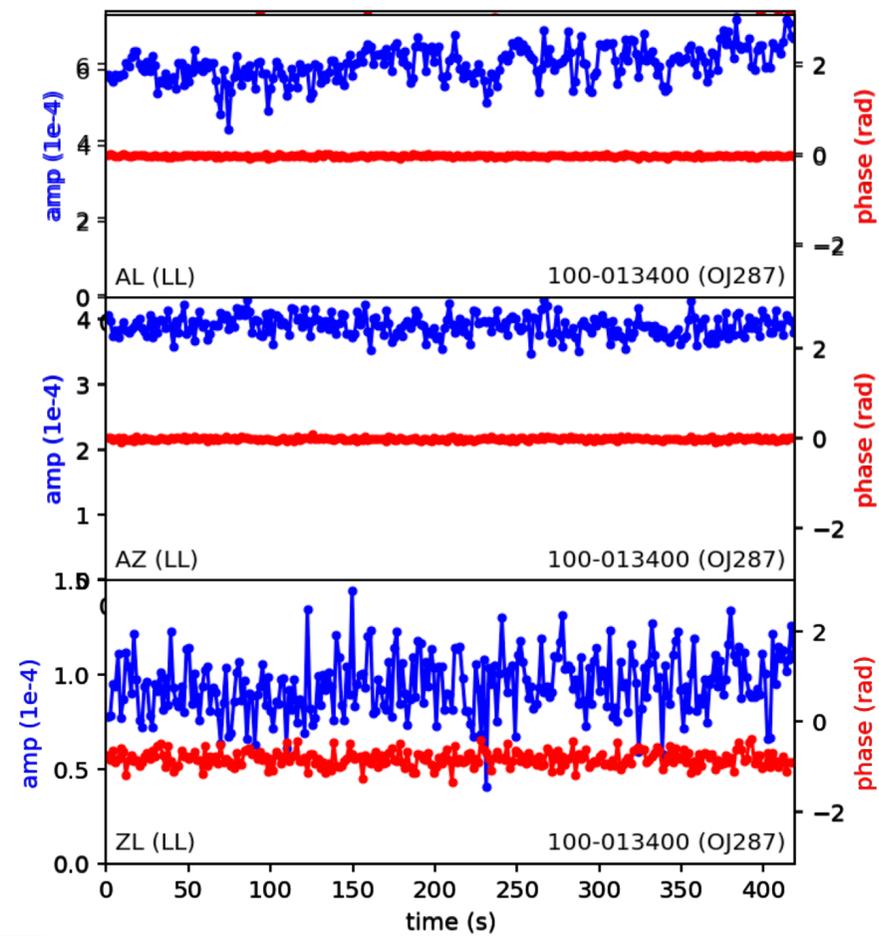
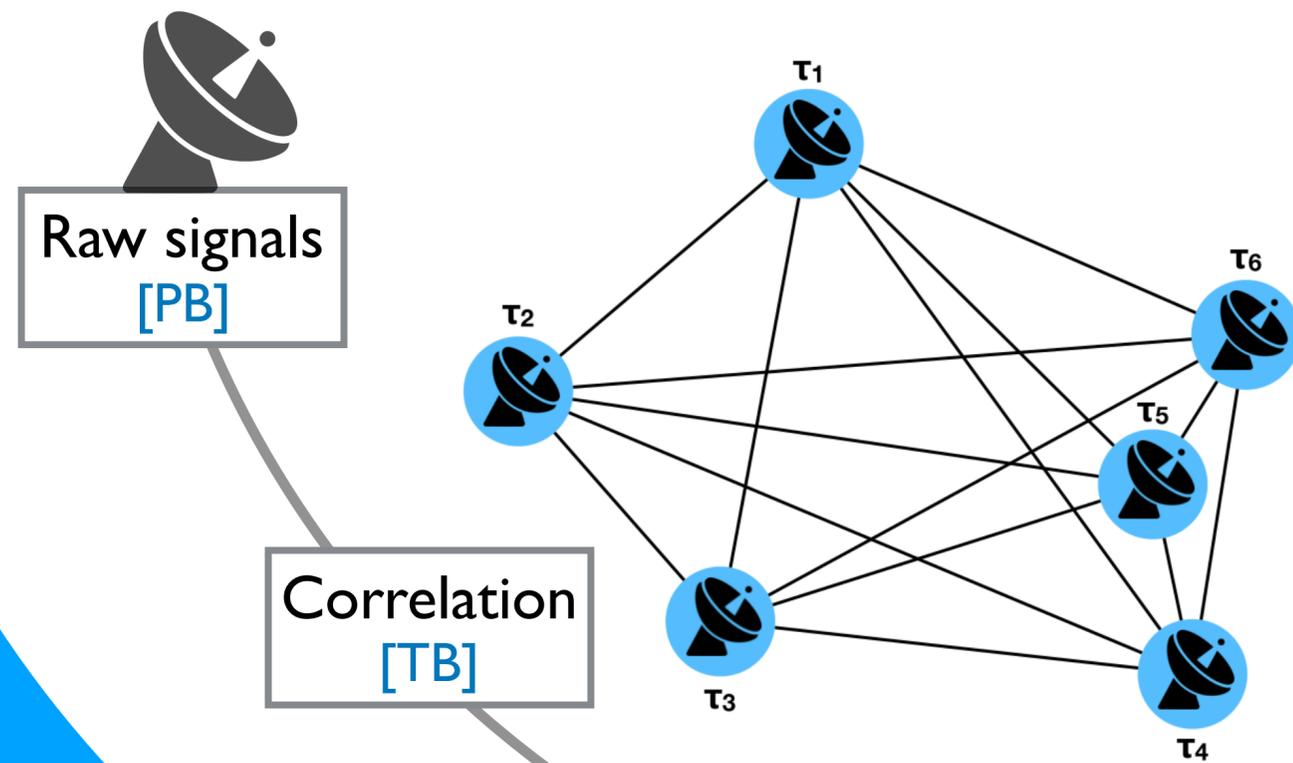
Reid & Honma

We just need to measure E_1 and E_2 at various locations in the plane of propagation, but..

1. Earth is round & moving
2. Irregular delays from troposphere/ionosphere
3. Different atmospheric and receiver noise
4. Various electronics and path delays
5. Independent and imperfect clocks at all stations
6. Post-digitization artifacts
7. Unexpected data issues

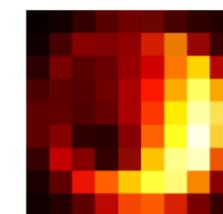
In data reduction, we are asked to “hide” as many of these effects as possible (without ruining the data)

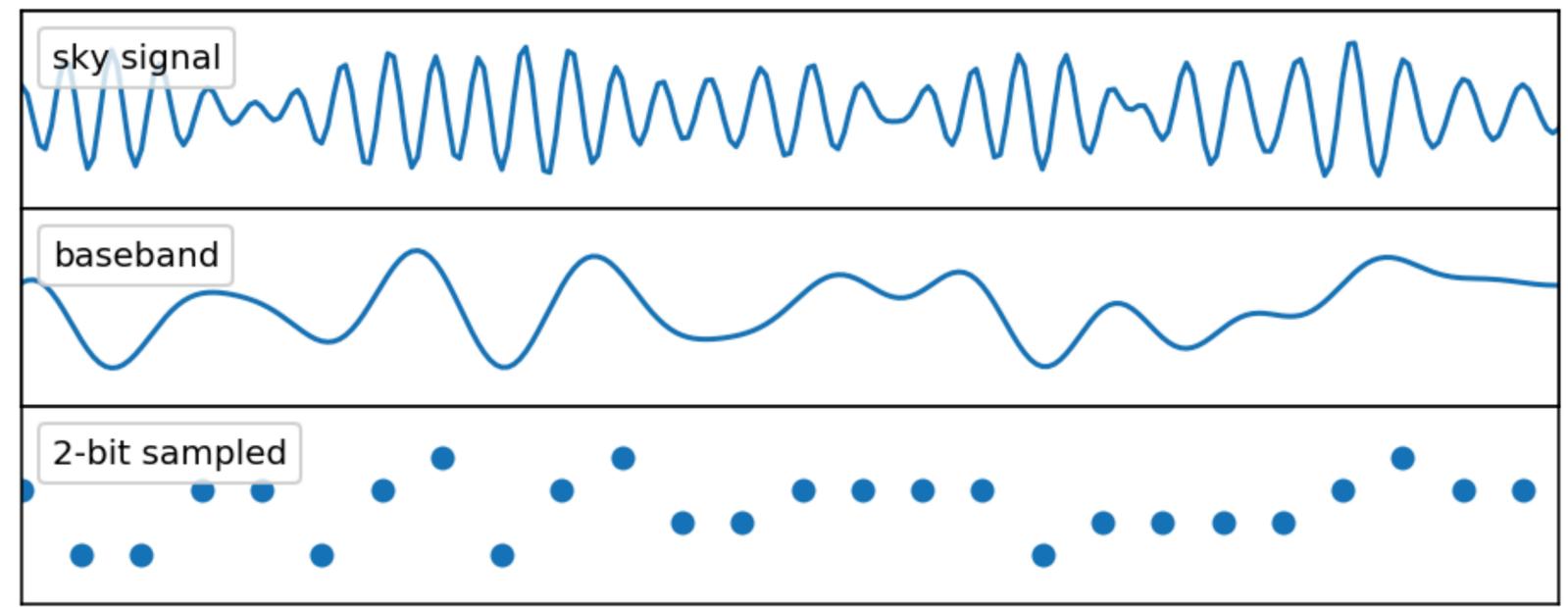
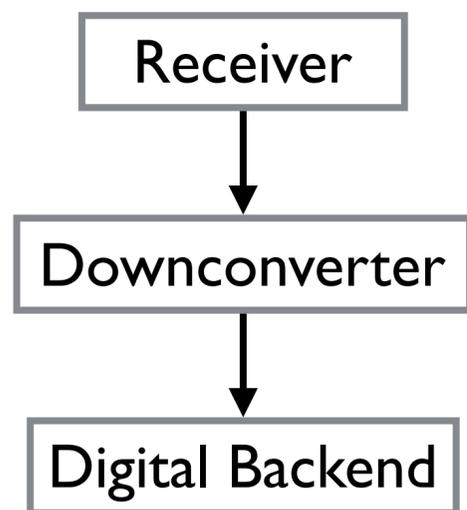
VLBI data and calibration pathway



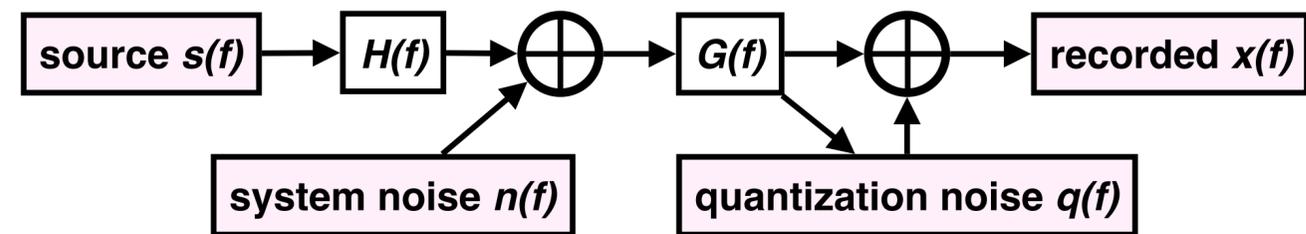
Calibration
[MB]

Analysis
[kB]





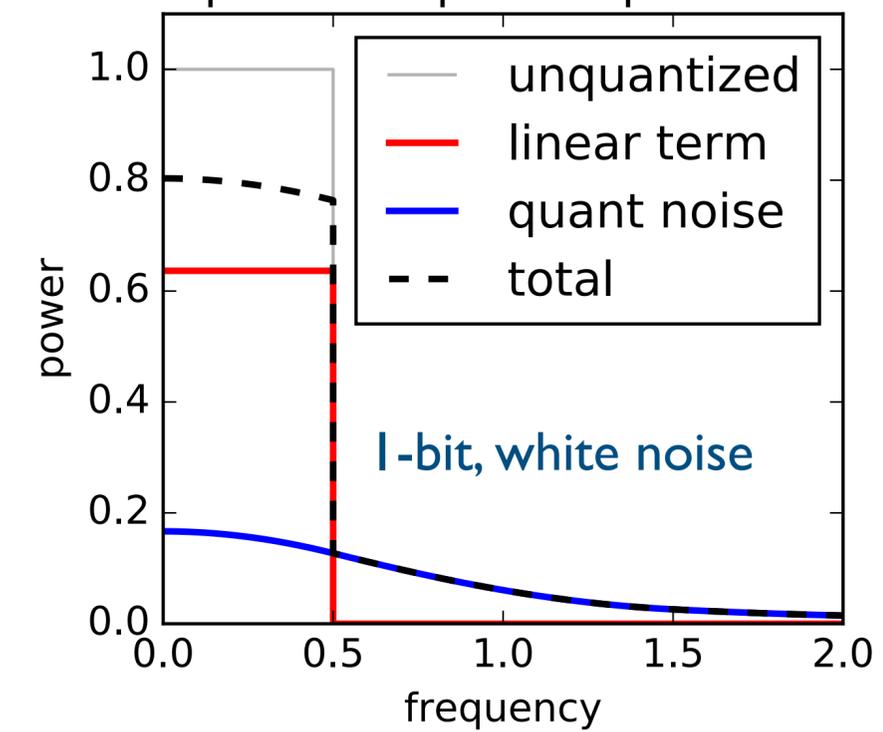
Linear components (bandpass, delay, dispersion)



There are **two** important bandpass effects $H(f)$ and $G(f)$, sometimes factored into a real (autocorr) and complex BP

Non-linear effects (delay-rate, atmospheric phase) must be described using **time-dependent** factors

quantized power spectrum





Gains, Polarization, and the Measurement Equation

Propagation of the astrophysical signal E through measurement v can be characterized by **complex** gain factors g

$$v(t, f) = g(t, f)E(t, f) \quad \langle v_1 v_2^* \rangle = g_1 g_2^* \langle E_1 E_2^* \rangle$$

Signal and ensemble averages are parameterized in **time** and **frequency**, which requires that g is varying (relatively) slowly

For two orthogonal feeds of an antenna, this can be written in matrix form,

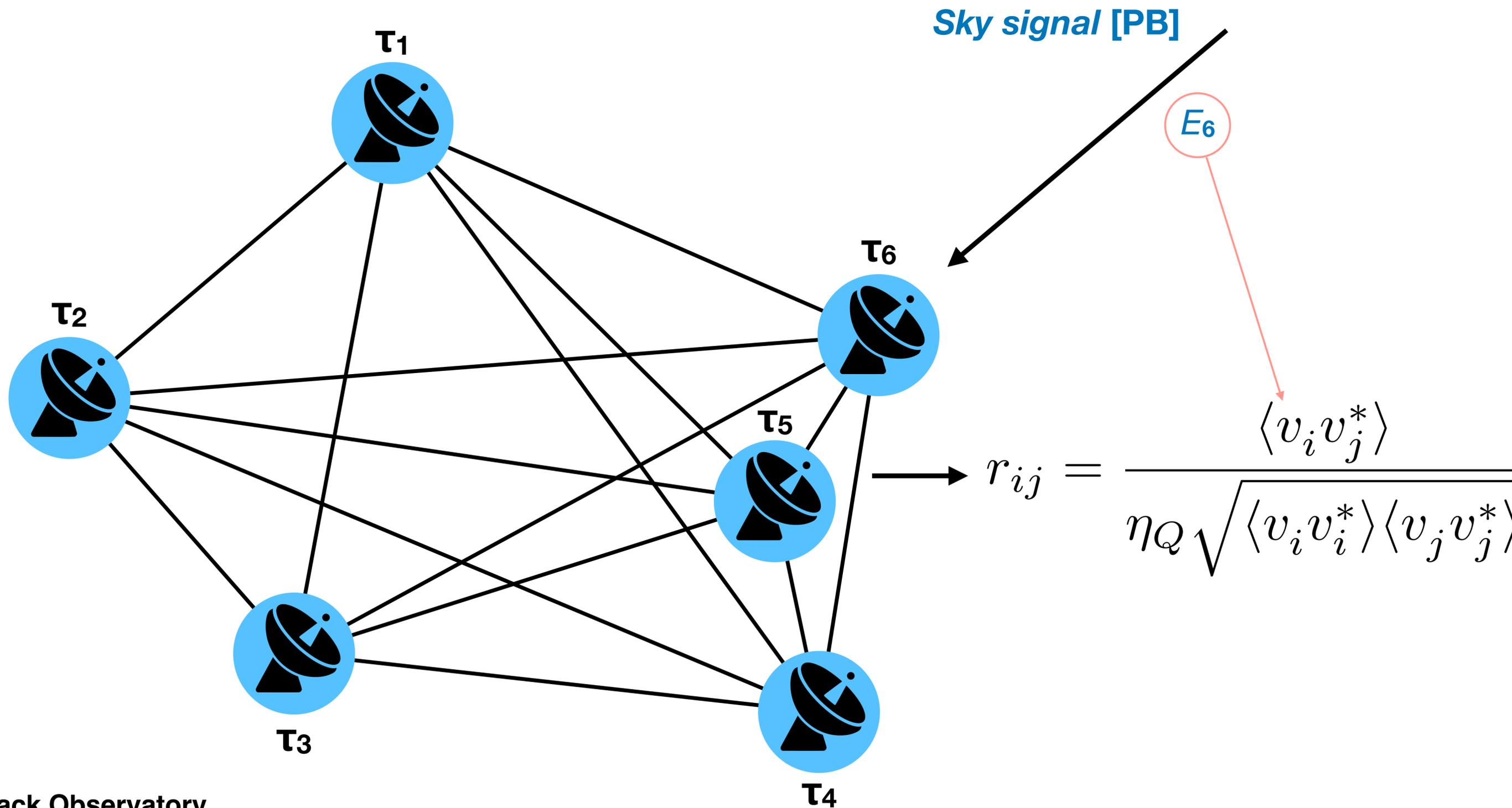
$$\begin{pmatrix} v_L \\ v_R \end{pmatrix} = \begin{pmatrix} g_L & 0 \\ 0 & g_R \end{pmatrix} \begin{pmatrix} E_L \\ E_R \end{pmatrix} \quad \begin{pmatrix} \langle v_{1L} v_{2L}^* \rangle & \langle v_{1R} v_{2L}^* \rangle \\ \langle v_{1R} v_{2L}^* \rangle & \langle v_{1R} v_{2R}^* \rangle \end{pmatrix} = \begin{pmatrix} g_{1L} & 0 \\ 0 & g_{1R} \end{pmatrix} \begin{pmatrix} \langle E_{1L} E_{2L}^* \rangle & \langle E_{1R} E_{2L}^* \rangle \\ \langle E_{1R} E_{2L}^* \rangle & \langle E_{1R} E_{2R}^* \rangle \end{pmatrix} \begin{pmatrix} g_{2L}^* & 0 \\ 0 & g_{2R}^* \end{pmatrix}$$

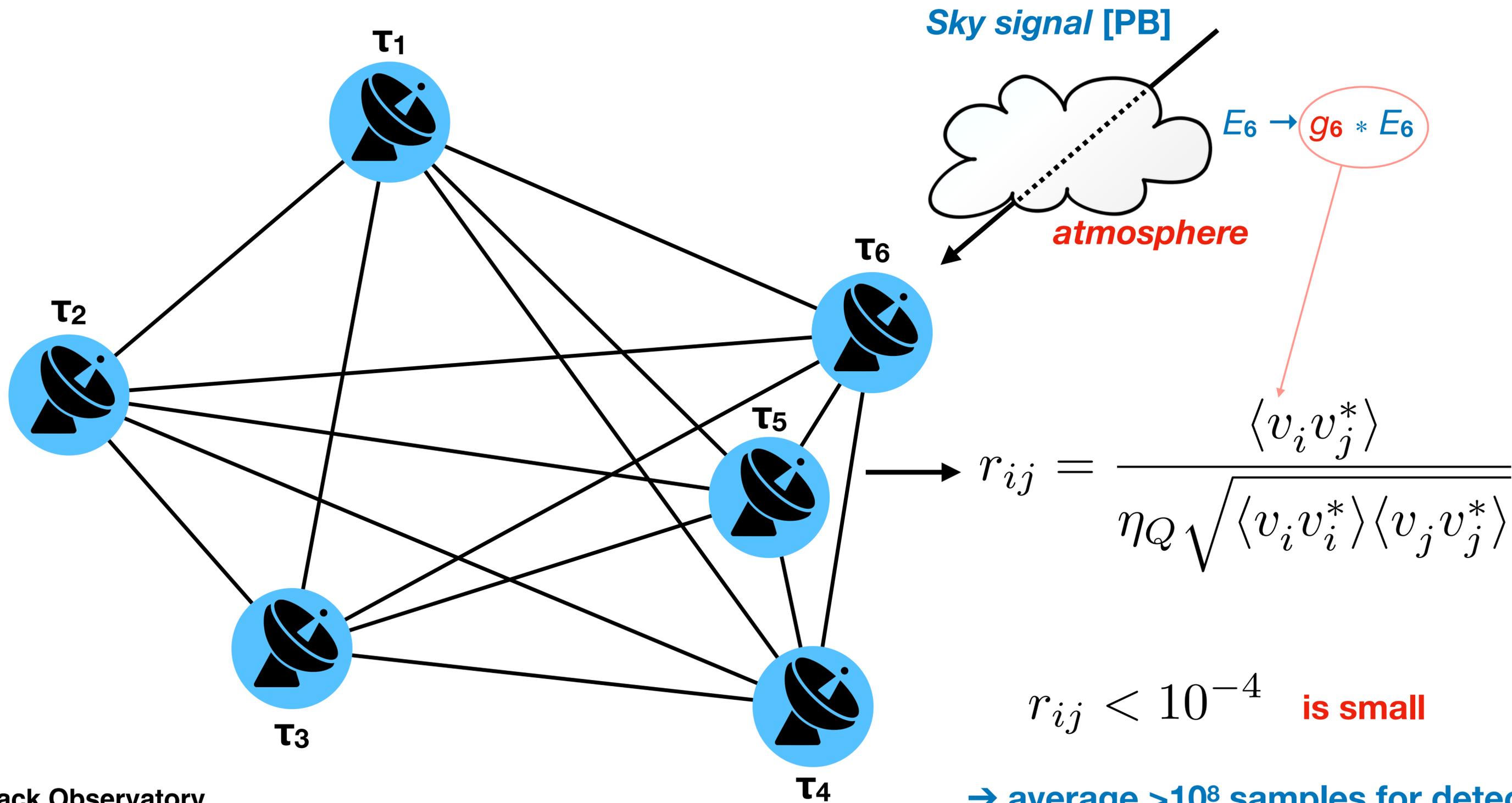
Tracking various physical propagation effects, as well as non-zero off-diagonal “D” terms (leakage across feeds, or change of polarization basis), leads to **Jones matrix** formalism used by the **Measurement Equation**

$$\mathbf{v} = \mathbf{J}_a \mathbf{J}_b \cdots \mathbf{J}_z \mathbf{E} \quad \langle \mathbf{v}_1 \mathbf{v}_2^\dagger \rangle = \mathbf{J}_{1a} \mathbf{J}_{1b} \cdots \mathbf{J}_{1z} \langle \mathbf{E}_1 \mathbf{E}_2^\dagger \rangle \mathbf{J}_{2z}^\dagger \cdots \mathbf{J}_{2a}^\dagger \mathbf{J}_{2b}^\dagger \quad (\text{see } \text{Smirnov 2011})$$

Why so many? **Physical model** generally allows for least complexity. Note that matrices do not necessarily commute!

This is a very useful structure! One still must adopt good **models** for all the Jones matrices.. also track **noise**..







Flux calibration (a priori)

The **correlation coefficient** is normalized by the system noise in the separate receiving systems

Relating this to physical units of **correlated flux density** requires a calibration of the noise power

$$r_{ij} = \frac{\langle v_i v_j^* \rangle}{\eta_Q \sqrt{\langle v_i v_i^* \rangle \langle v_j v_j^* \rangle}} \quad |V_{ij}| = \sqrt{\text{SEFD}_i \times \text{SEFD}_j} |r_{ij}|.$$

This is encapsulated into the **system-equivalent flux density** (SEFD) at each site, which is the (measured) noise power in units of flux density from an unpolarized astrophysical source (above the atmosphere)

The SEFD is calibrated separately from the data using **first principles**, known bright calibrators (**planets**), and **noise sources** of known temperature placed directly in front of receiving elements, and is taken “a priori”

For a heterogeneous array such as the EHT, SEFD can range by **orders of magnitude** $\sim 10^2$ to $\sim 10^5$ Jy



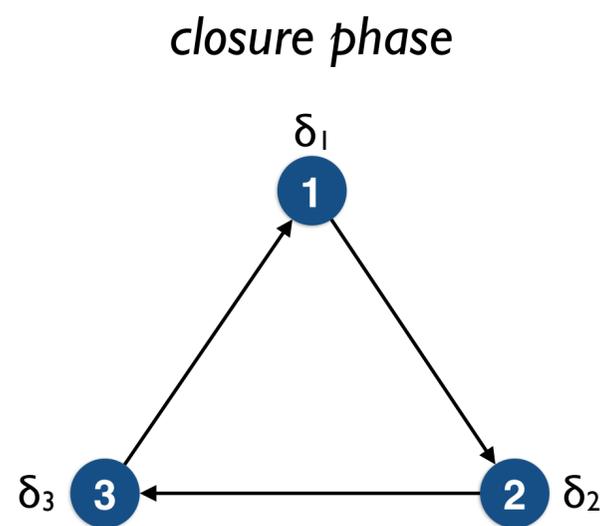
Closure relationships

At **mm-frequencies**, phase transfer from nearby calibration targets is very difficult or impossible so we have essentially **no a priori information** about station phase

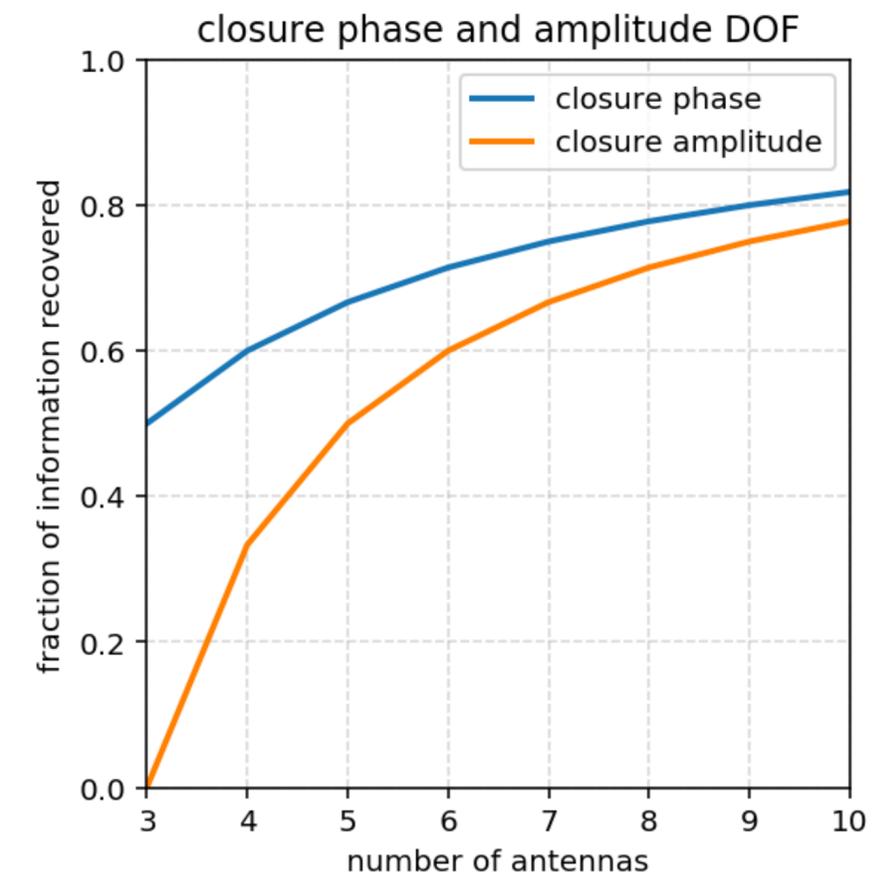
$$r_{12} = \frac{\langle x_1 x_2^* \rangle}{\eta_Q \sqrt{\langle x_1 x_1^* \rangle \langle x_2 x_2^* \rangle}} = \frac{e^{i\theta_1} e^{-i\theta_2} \mathcal{V}_{12}}{\sqrt{\text{SEFD}_1 \times \text{SEFD}_2}}$$

However there are $N(N-1)/2$ baseline measurements of phase, yet only $(N-1)$ unknown station phases, so the measurements do capture structural phase information about the source

This information is captured by the “**closure phases**”



*insensitive to relative phase of each antenna:
N-1 degrees of freedom removed from baselines*

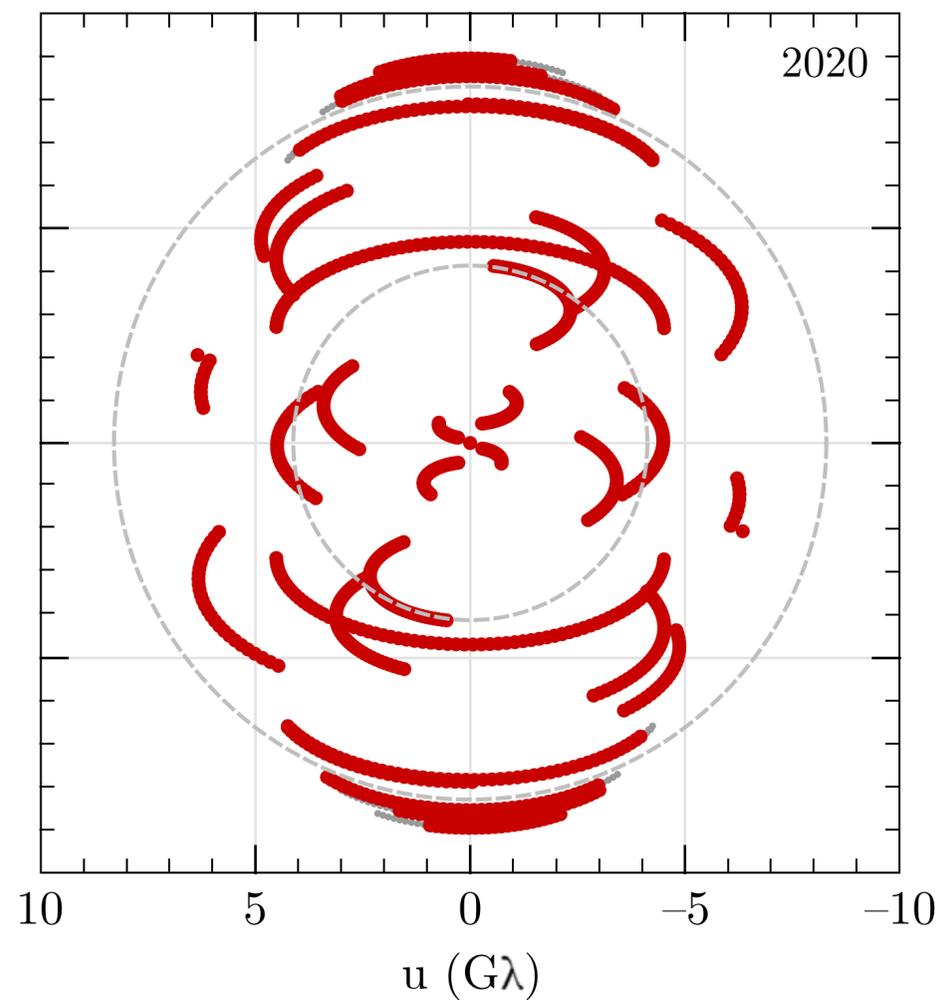




Time-frequency ensemble average: source constraints

The correlation coefficients measured by the interferometer relies on finite averages to estimate expectation value

$$\langle v_1 v_2^* \rangle = g_1 g_2^* \langle E_1 E_2^* \rangle$$



EHTC 2019 ApJL 875 (Paper II)

What are the limits from the source?

$$\langle E_1 E_2^* \rangle = \iint e^{-2\pi \mathbf{u} \cdot \boldsymbol{\sigma}} I_\nu(\boldsymbol{\sigma}) d\Omega$$

$$= \mathcal{V}(\mathbf{u})$$

Interferometer sweeps through \sim FOV/beam measurements in 24h
For EHT sources of \sim few² independent pixels, coherence length \sim hours

A \sim few pixels across a spatial dimension means $>10\%$ fractional bandwidth can be averaged without affecting independent measurements

Compact EHT sources implies **intrinsic smoothness/stability**
in time and frequency for the model visibility

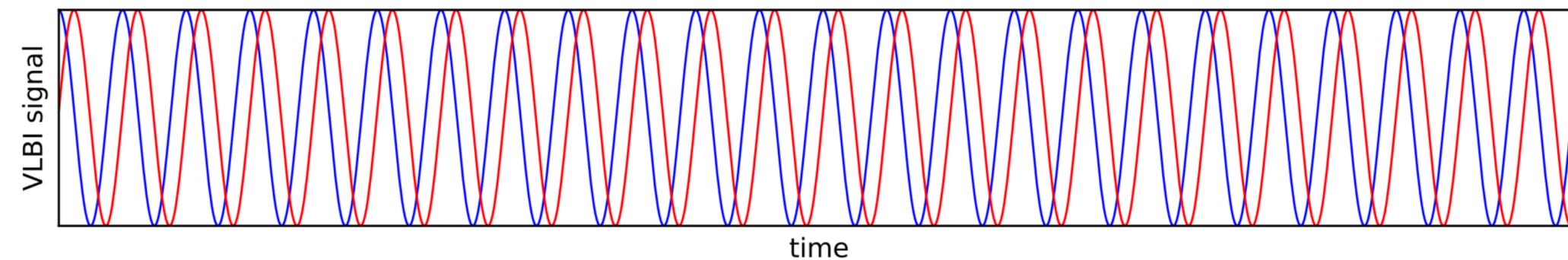
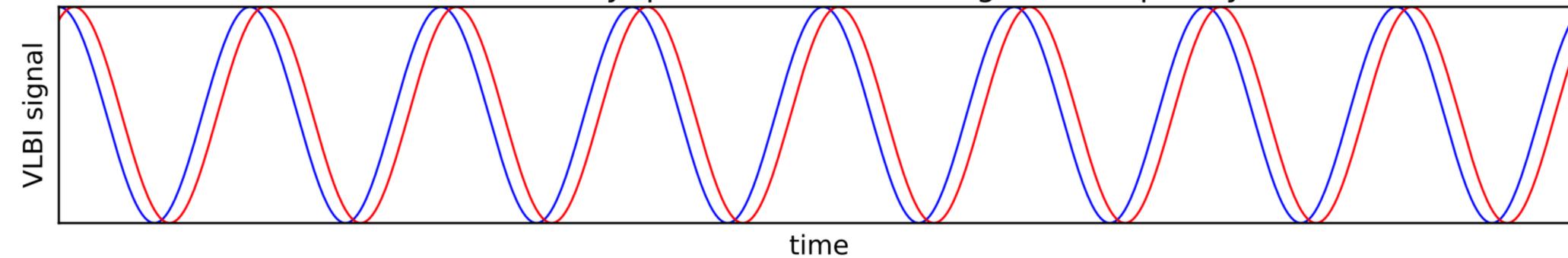


Time-frequency ensemble average: phase systematics

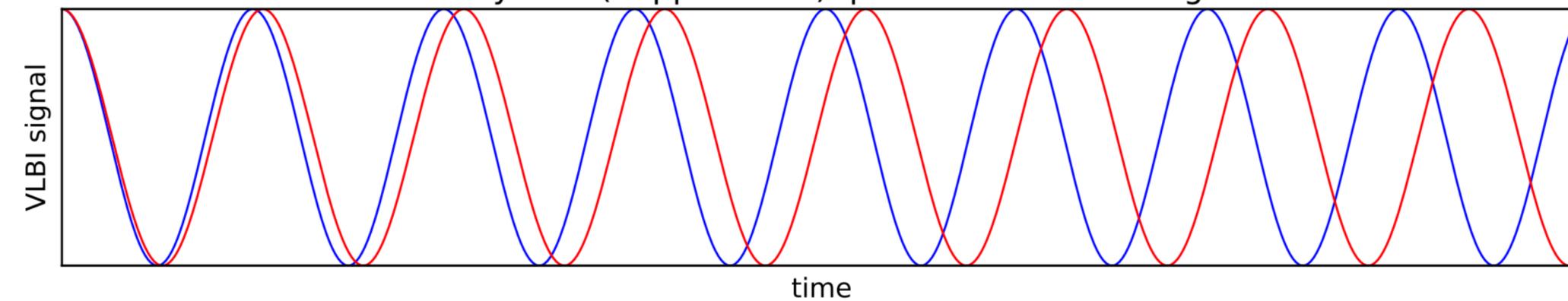
What about variability in gain parameters?

$$\langle v_1 v_2^* \rangle = g_1 g_2^* \langle E_1 E_2^* \rangle$$

effect of delay: phase shift increasing with frequency



effect of delay-rate (doppler shift): phase shift increasing with time



First-order phase systematics

$$\Delta\phi = \frac{\partial\phi}{\partial\nu} \Delta\nu + \frac{\partial\phi}{\partial t} \Delta t$$

Delay Delay-rate (rate)

$$\tau = \frac{1}{2\pi} \frac{\partial\phi}{\partial\nu} \quad \dot{\tau} = \frac{1}{2\pi\nu} \frac{\partial\phi}{\partial t}$$

Large delays and rates taken out at Correlator at high time-frequency resolution using a **a priori Earth model** (calc)

We only worry about **residual clock errors**

0.5s × 0.5 MHz dump time, bandwidth

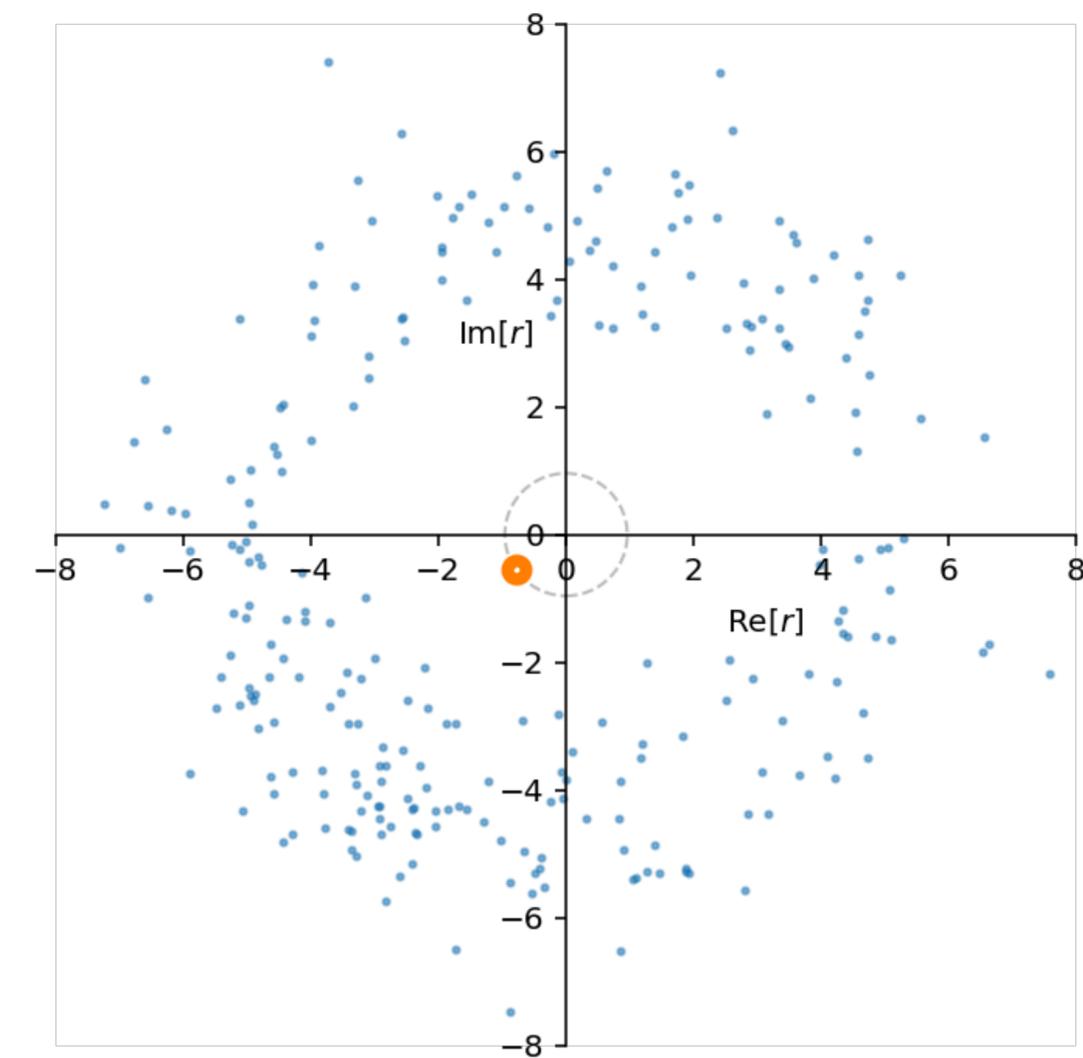
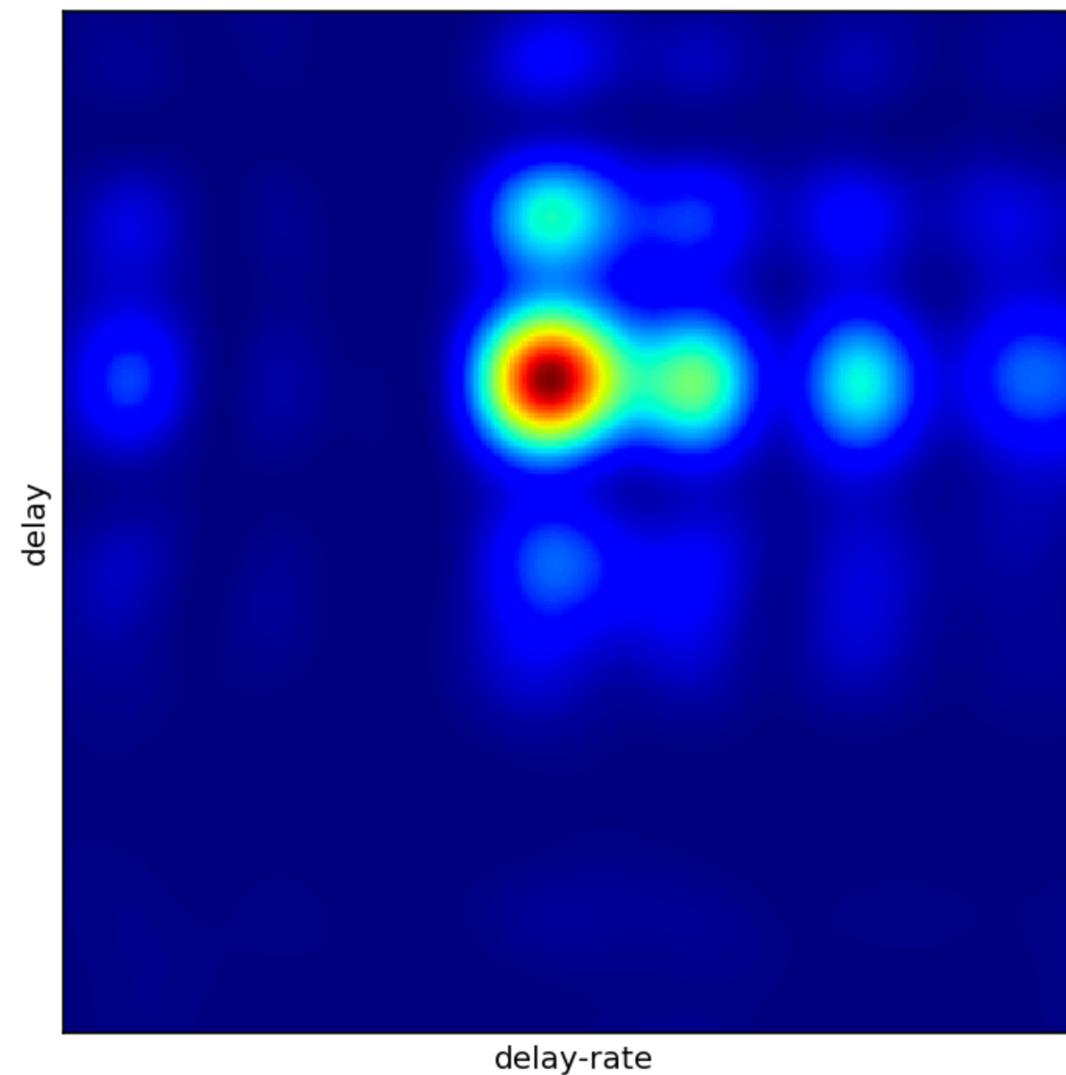
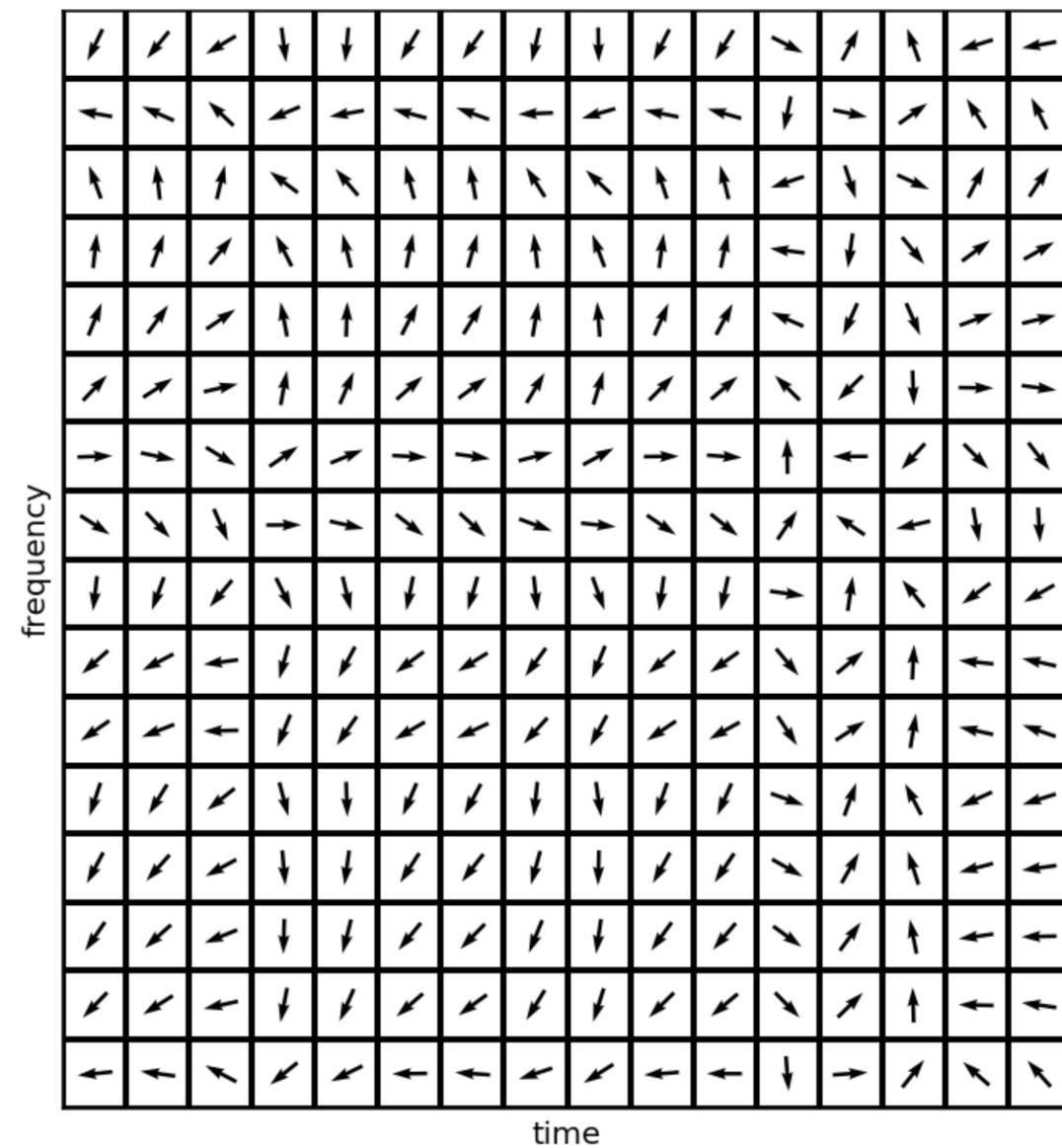
→ rates within ~2 ps/s (1.3 mm/s)

→ delays within ~1 μs

Fringe fitting

Fringe fitting involves **self-calibration of residual clock errors** to extract and average **correlation coefficient**
At high frequencies, there are linear and non-linear residuals in **phase vs frequency** and **phase vs time**

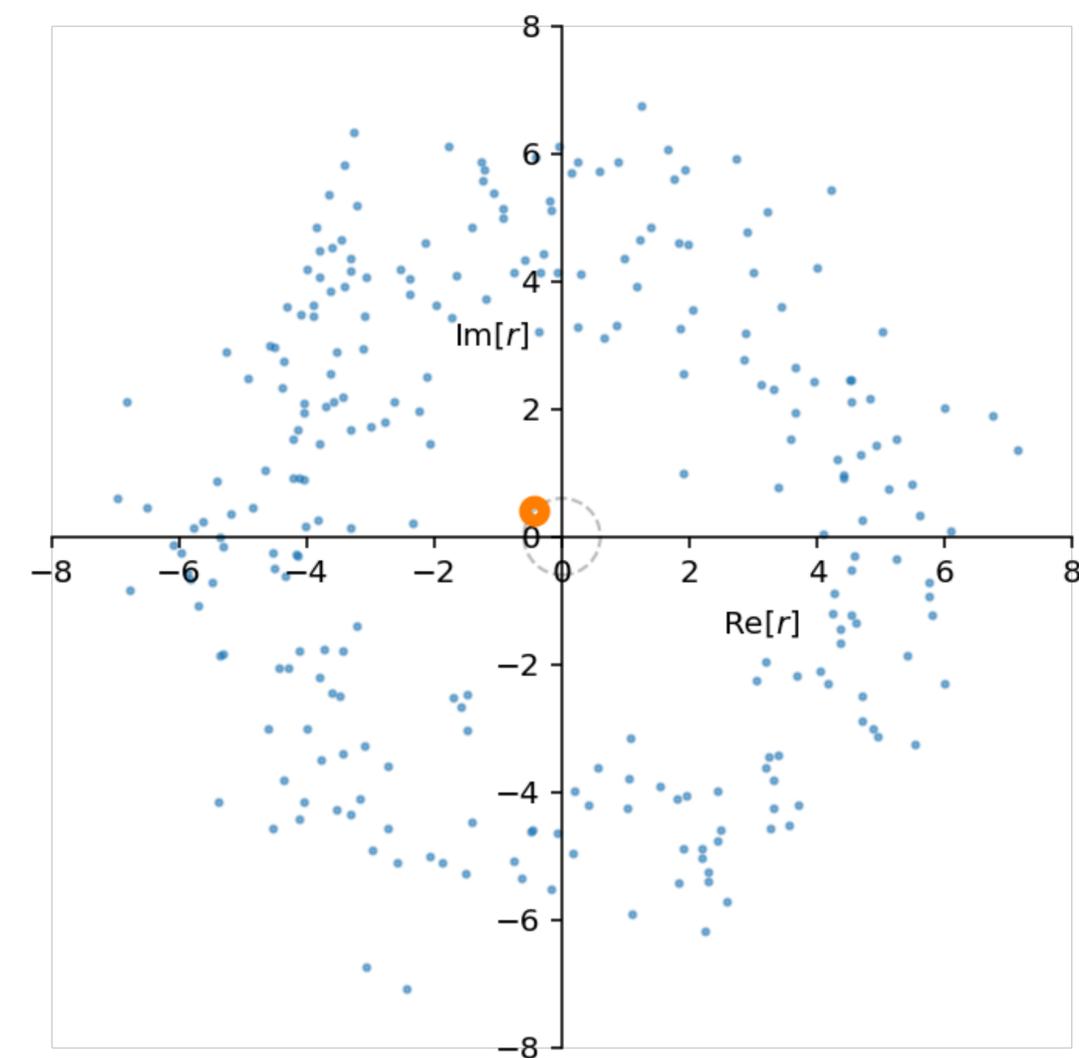
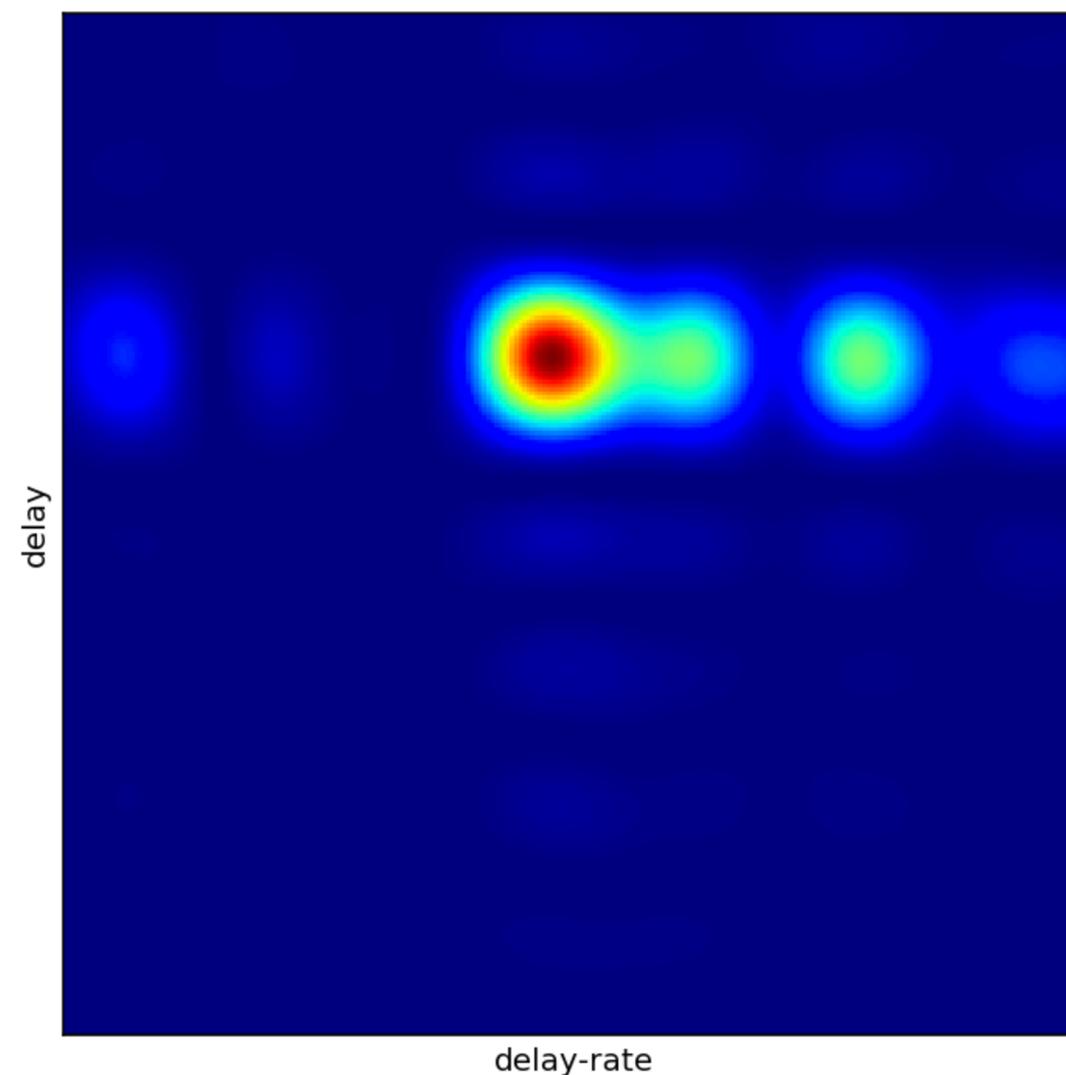
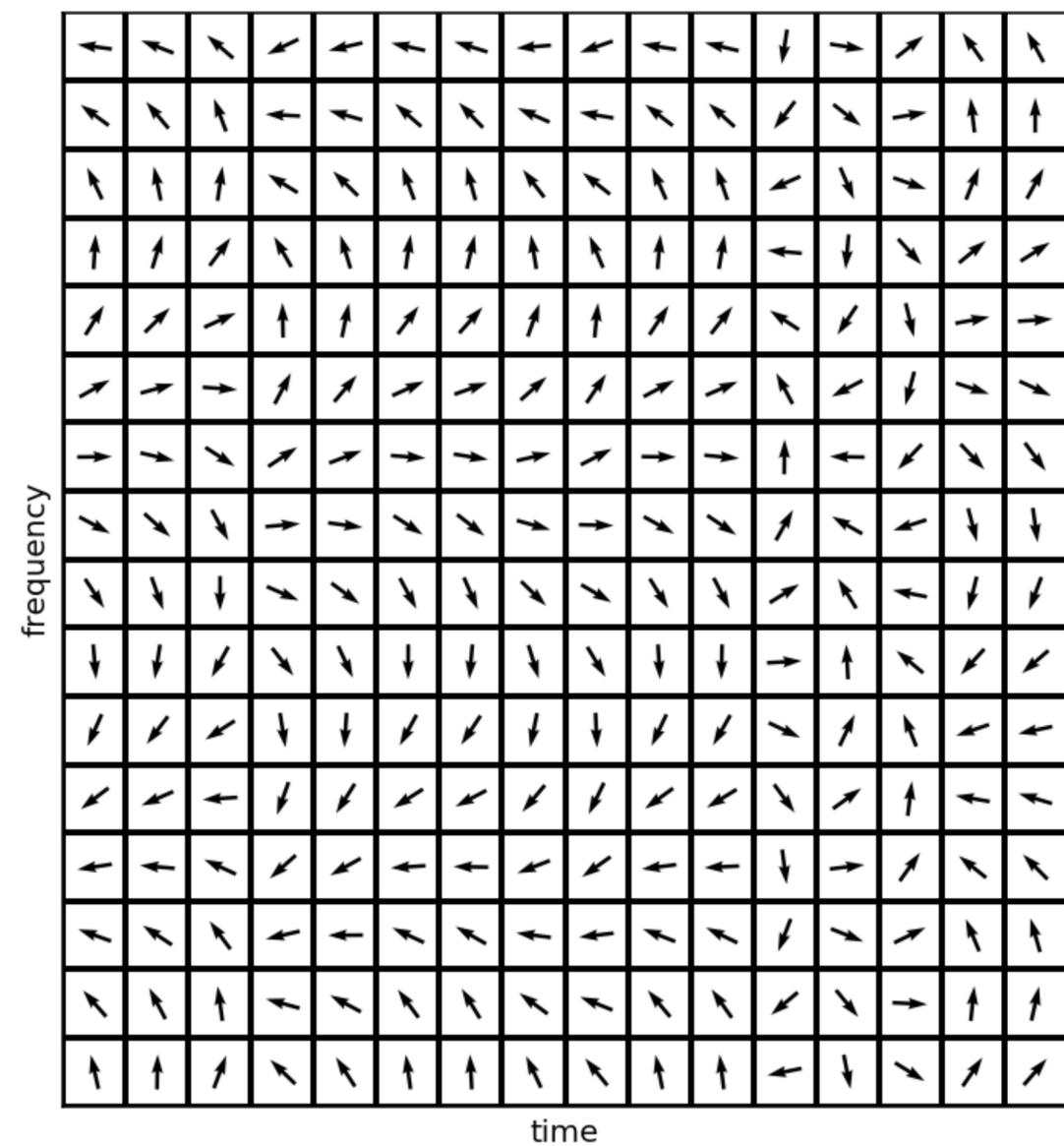
$$\Delta\phi_{12}(t, f, pp) = \phi_0 \quad (\text{a priori phase corrections})$$



Fringe fitting: phase bandpass

First correction is generally an **instrumental phase** bandpass because
It is **stable across the experiment** and can be solved on an ensemble of bright calibrators

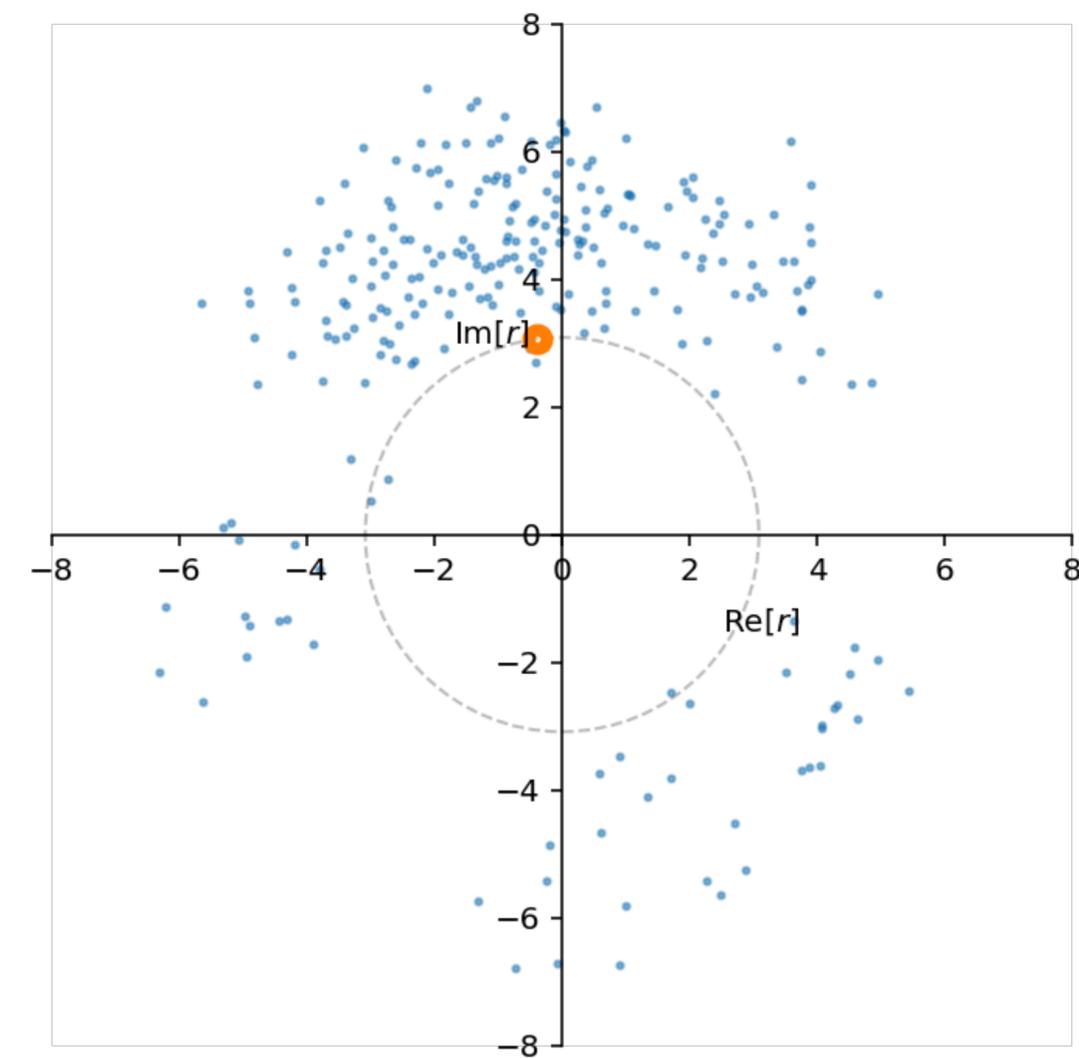
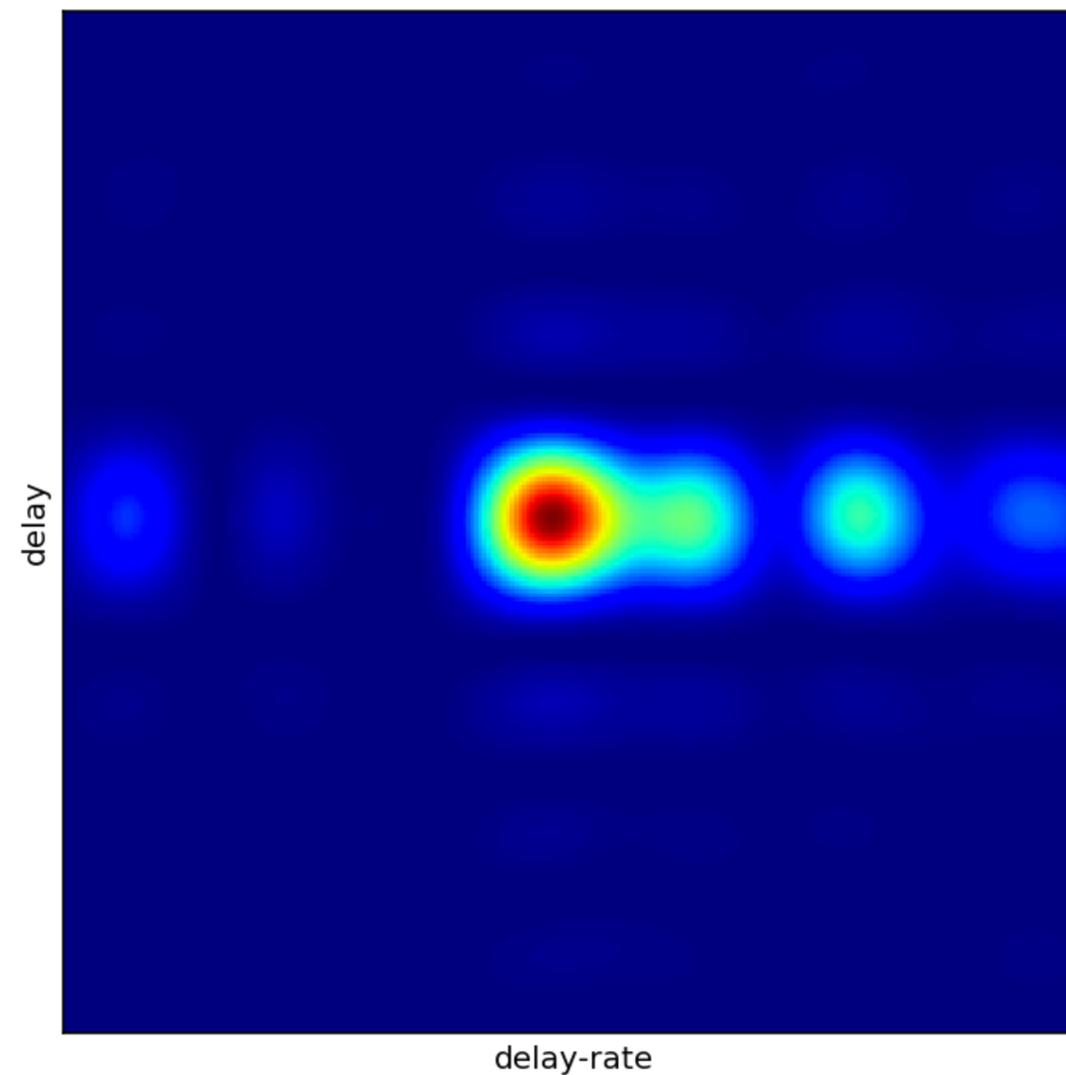
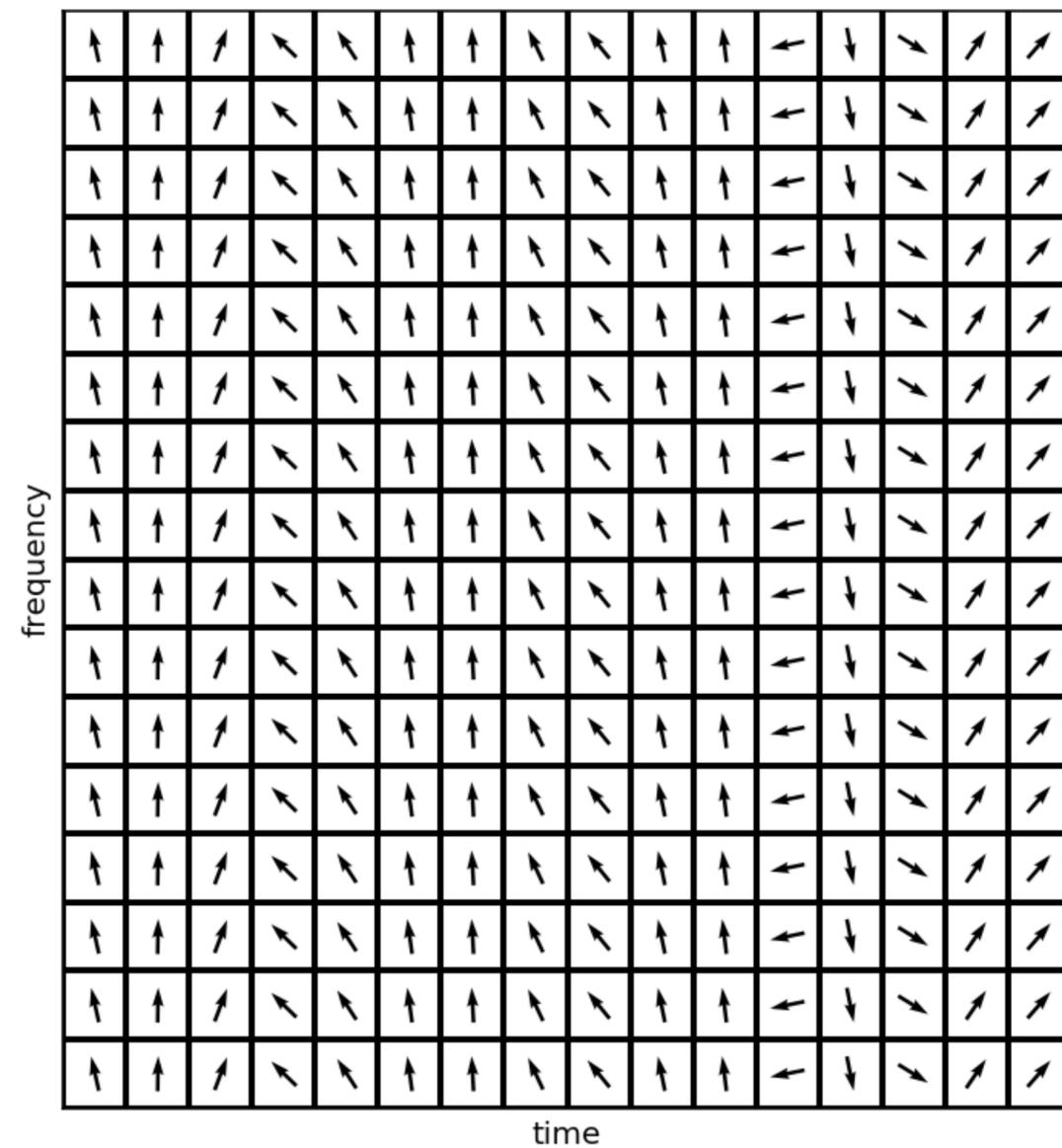
$$\Delta\phi_{12}(t, f, pp) = \phi_0 + \phi_{2-1}(f)$$



Fringe fitting: delay

After removing non-linear phase vs frequency, we can extract a clean linear fit to **delay** for this scan

$$\Delta\phi_{12}(t, f, \text{pp}) = \phi_0 + \phi_{2-1}(f) + 2\pi(f-f_{\text{ref}})\tau_{,\text{pp}}$$

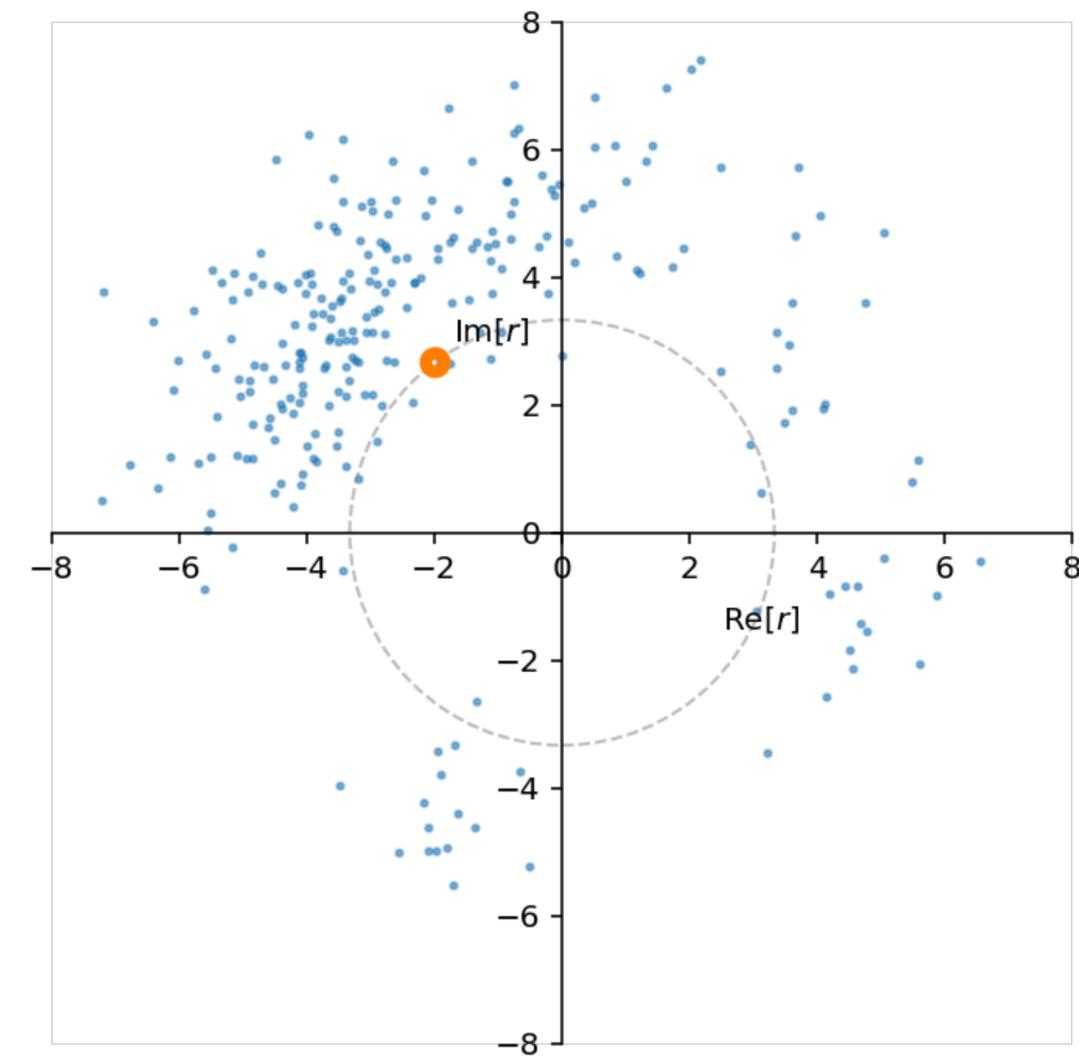
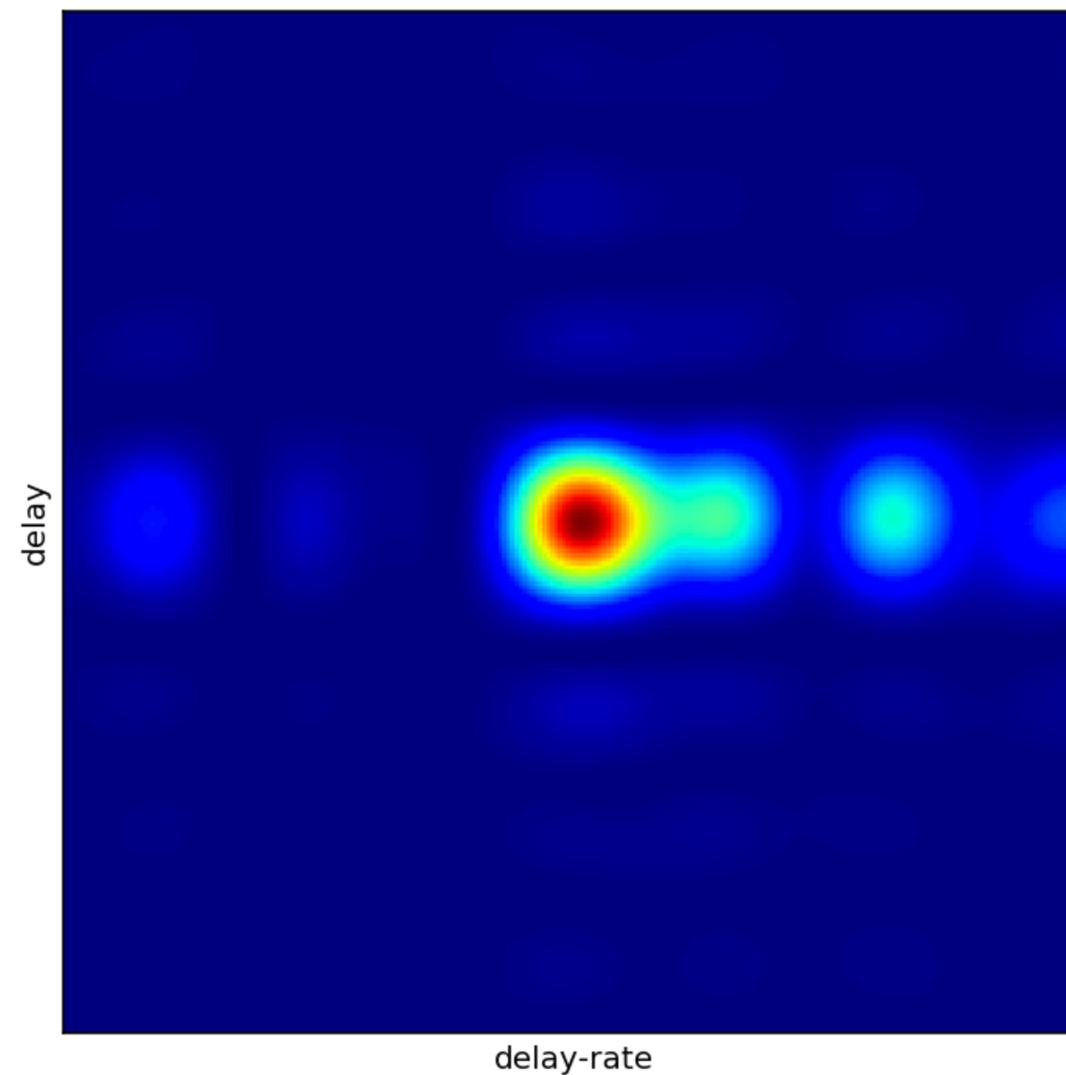
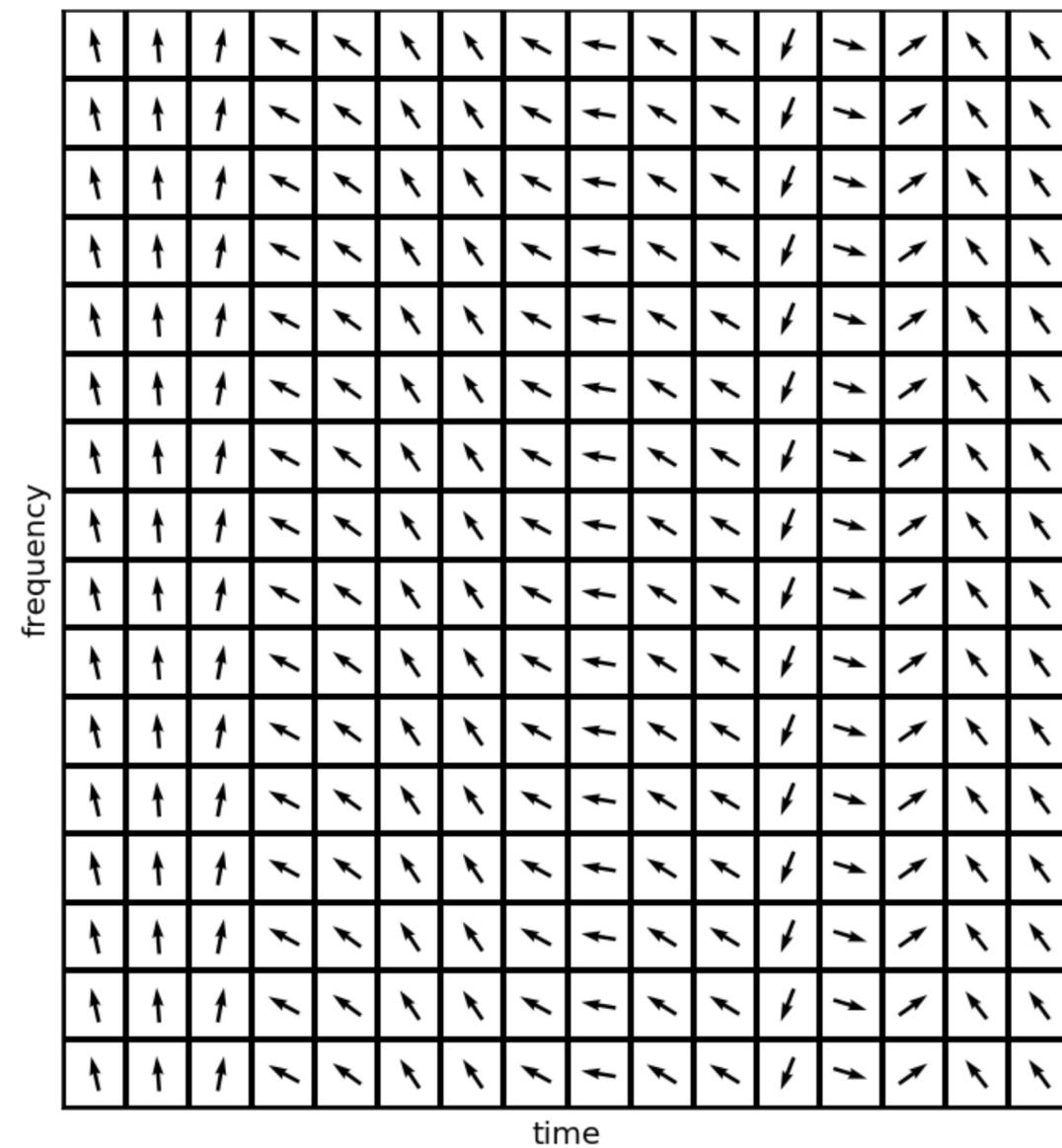




Fringe fitting: delay-rate

As well as **delay-rate**, although this is poorly defined in the presence of rapid atmospheric fluctuations

$$\Delta\phi_{12}(t, f, \text{pp}) = \phi_0 + \phi_{2-1}(f) + 2\pi(f-f_{\text{ref}})\tau_{\text{pp}} + 2\pi f(t-t_{\text{ref}})\dot{\tau}_{\text{pp}}$$

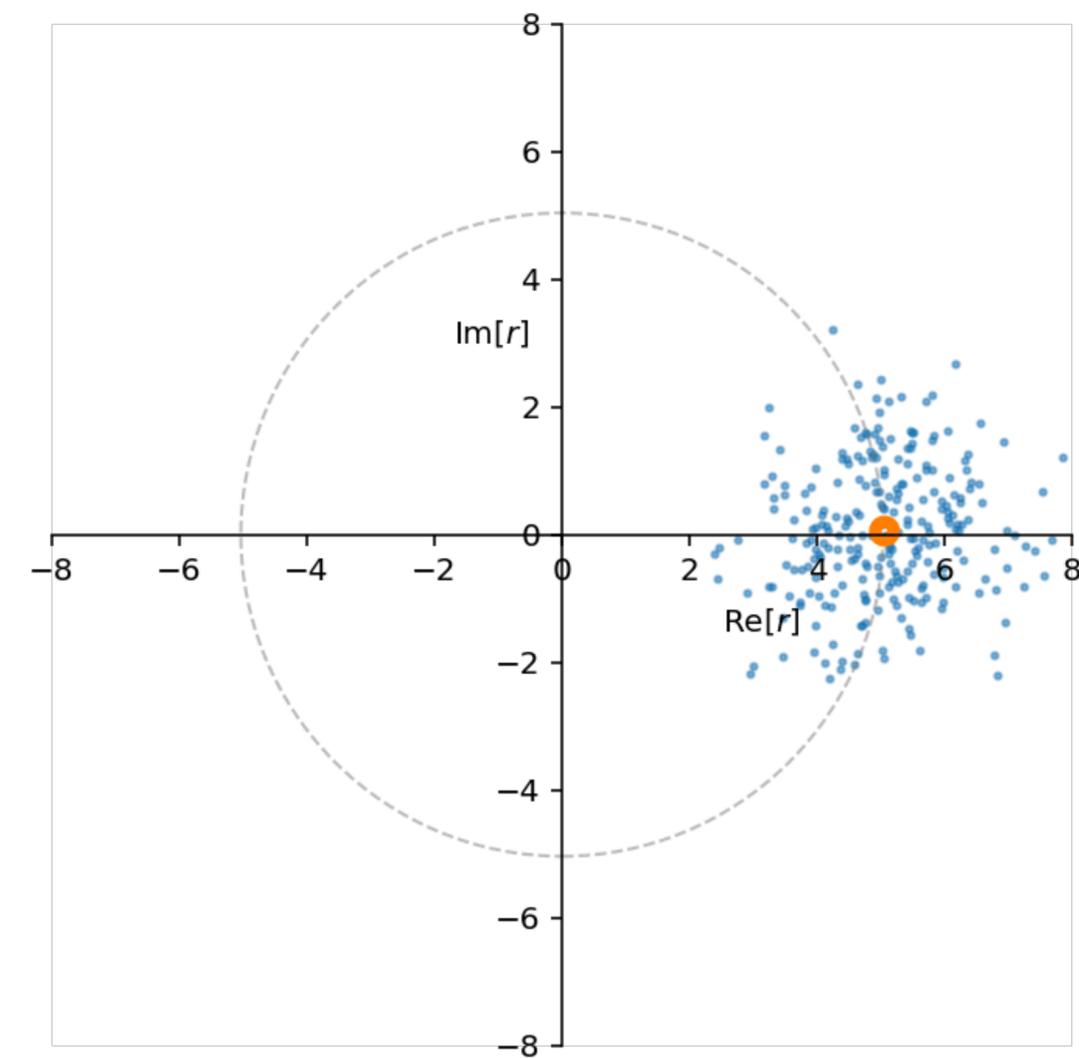
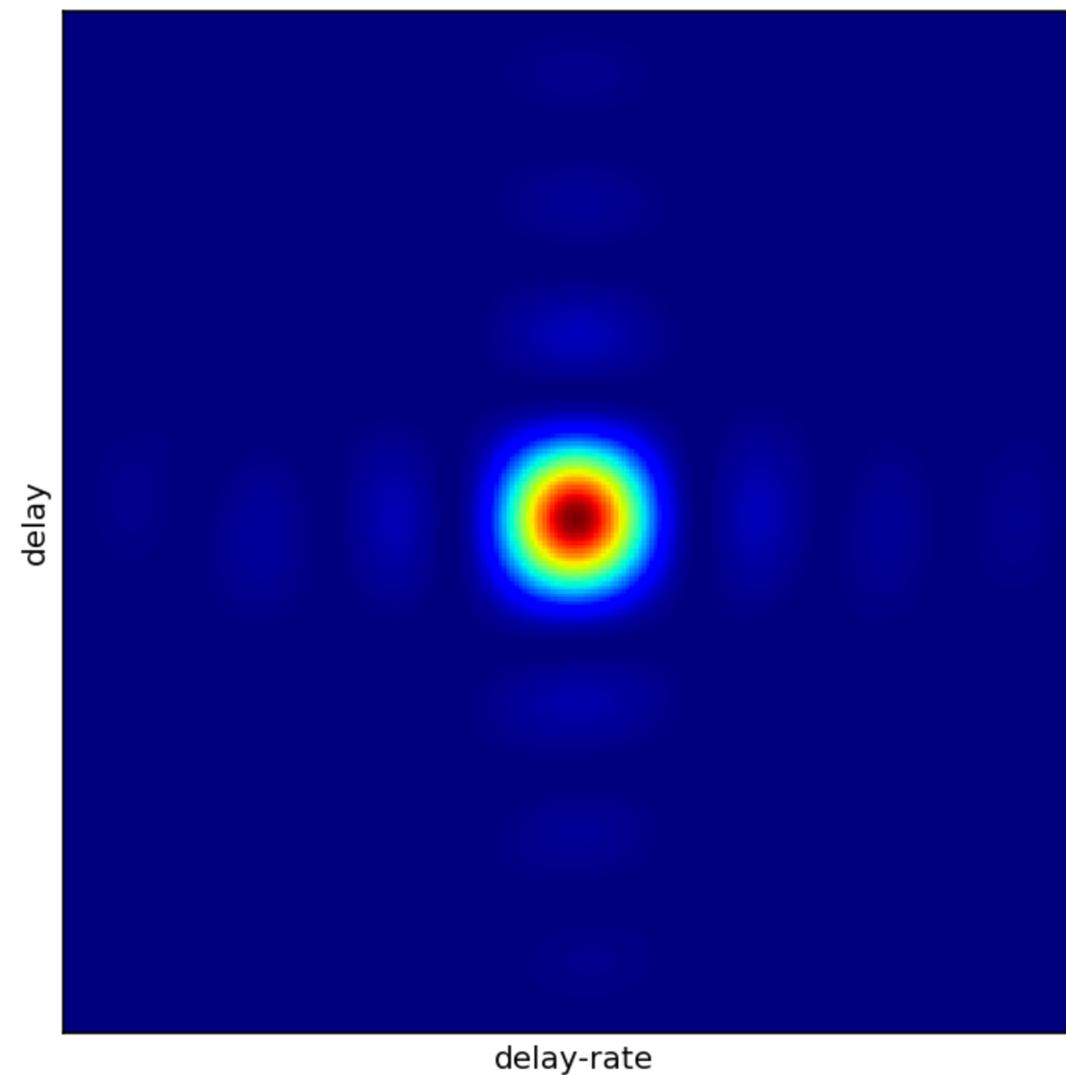
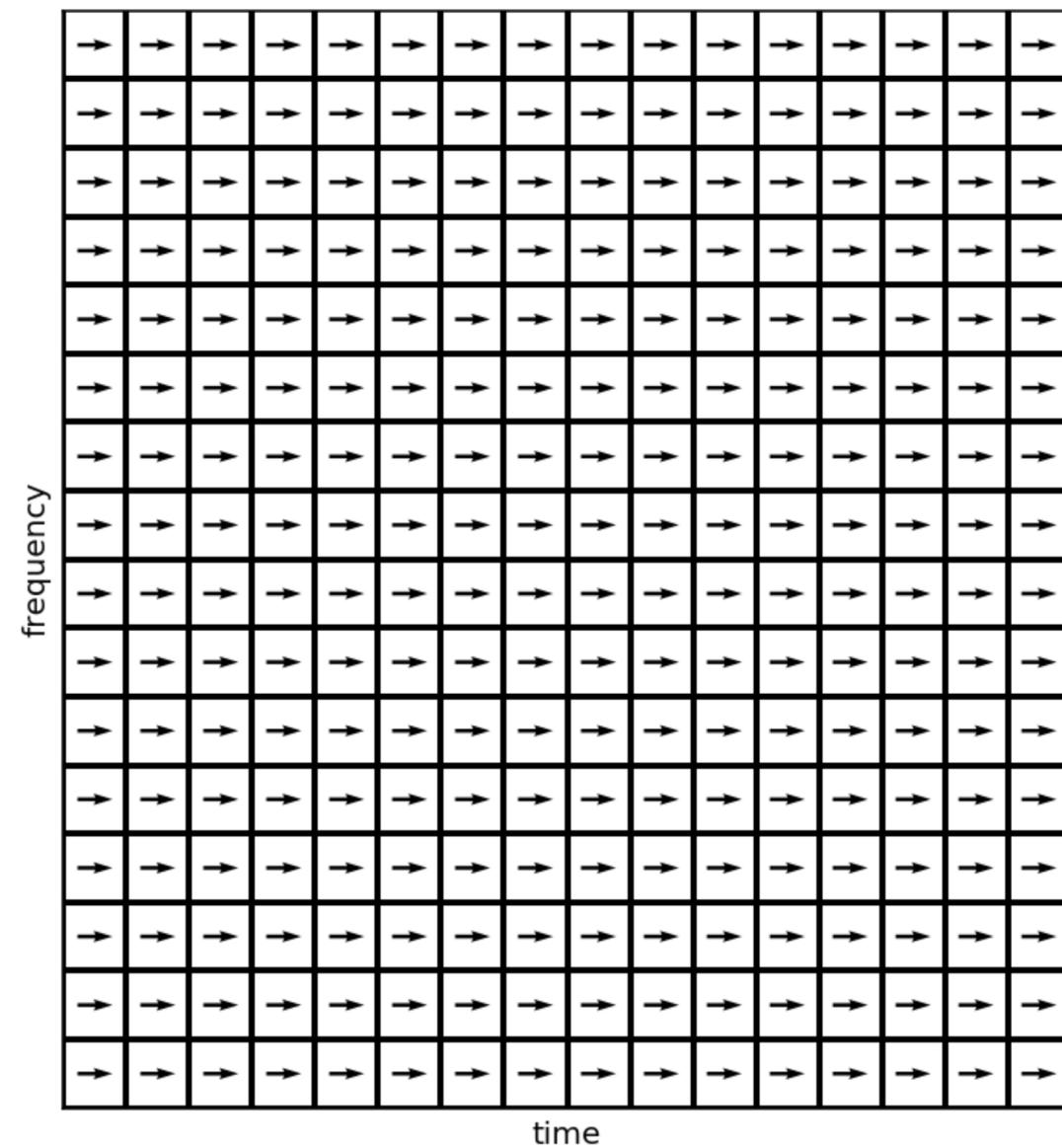


Fringe-fitting: atmospheric phase

And finally we can estimate and correct for **atmospheric phase**, here referencing to the first antenna

$$\Delta\phi_{12}(t, f, \text{pp}) = \phi_0 + \phi_{2-1}(f) + 2\pi(f-f_{\text{ref}})\tau_{\text{pp}} + 2\pi f(t-t_{\text{ref}})\dot{\tau}_{\text{pp}} + \phi_{2-1}(t)$$

now we can average over the entire scan and bandwidth

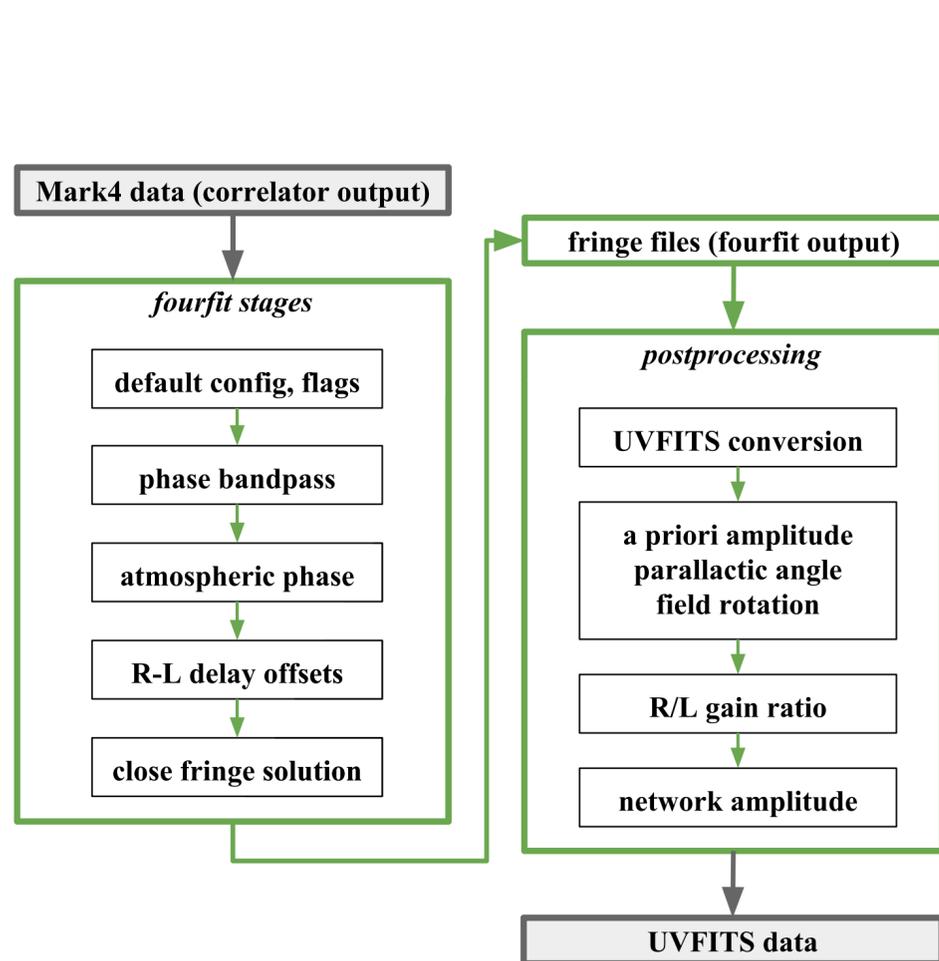




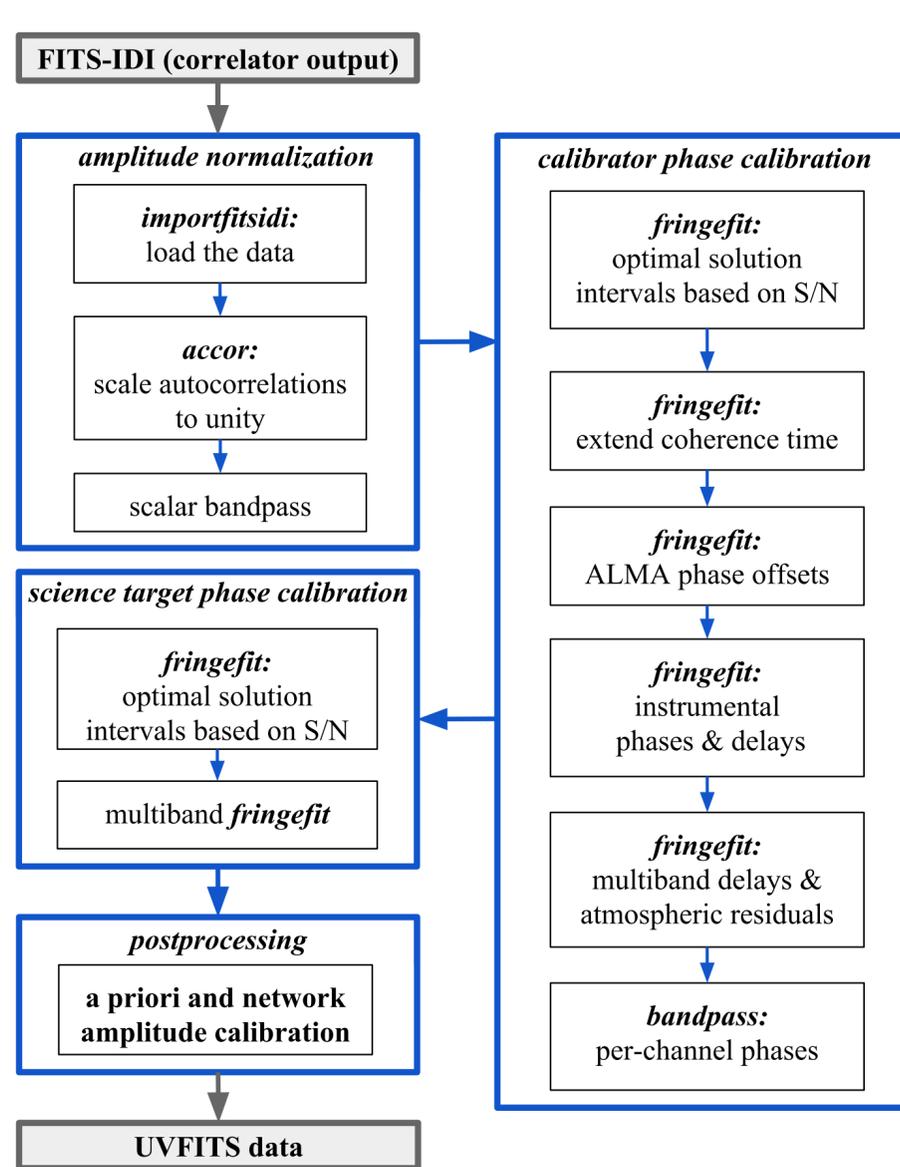
Phase calibration pipeline

For mm-VLBI such as EHT, **custom pipelines** are required due to uniqueness of data and systematics

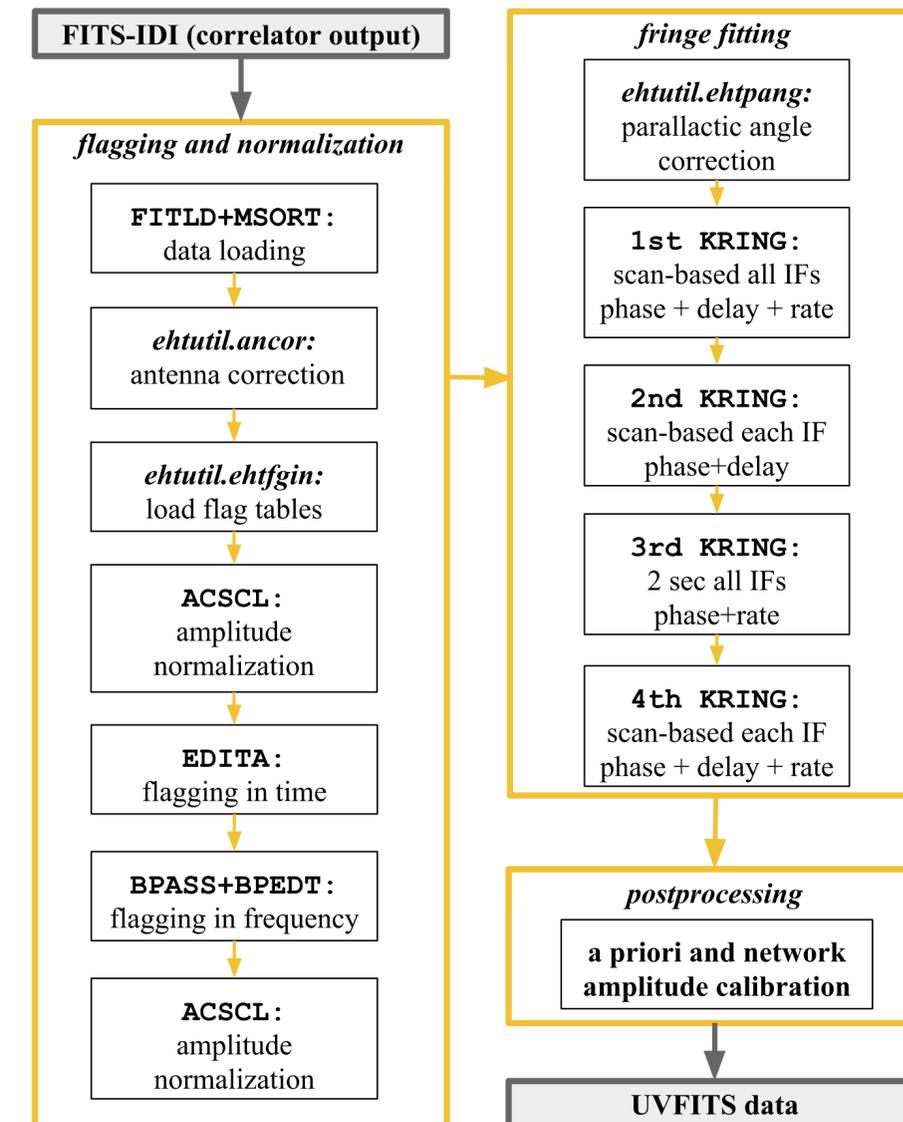
Purpose of steps is to fit as **simple** a model as possible, using as much **S/N** as available, and **maintain closure** (station-based gains)



“EHT-HOPS” (Blackburn+ 2019)



CASA “rPicard” (Janssen+ 2019)



AIPS (EHTC 2019 ApJL 875 (Paper III))

Some things that can go wrong

Too many free parameters for available S/N

Introduce calibration noise

Overfit data: bias amplitude upward, bias phase toward model

Underutilize array constraints and gain priors

Averaging over visibilities when gain is not stable

Introduce non-closing errors (averaged product of station gains may not factor)

Leaving in bad data / Ignoring systematics

Wrong calibration solutions

Systematic errors drive solution under the assumption of Gaussian thermal noise only

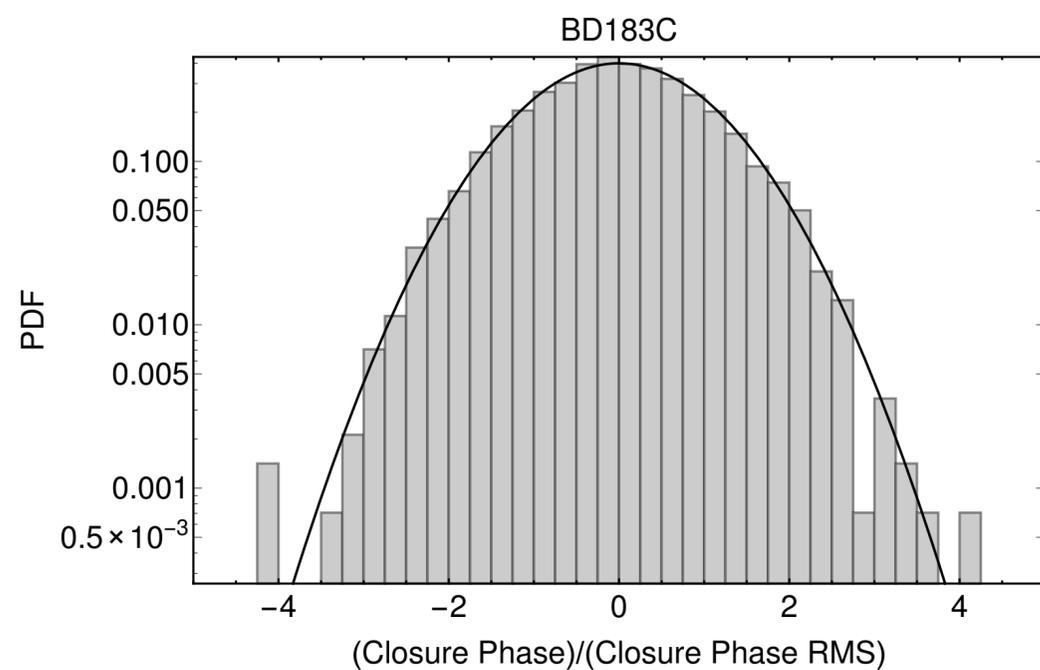


Thermal errors: origin

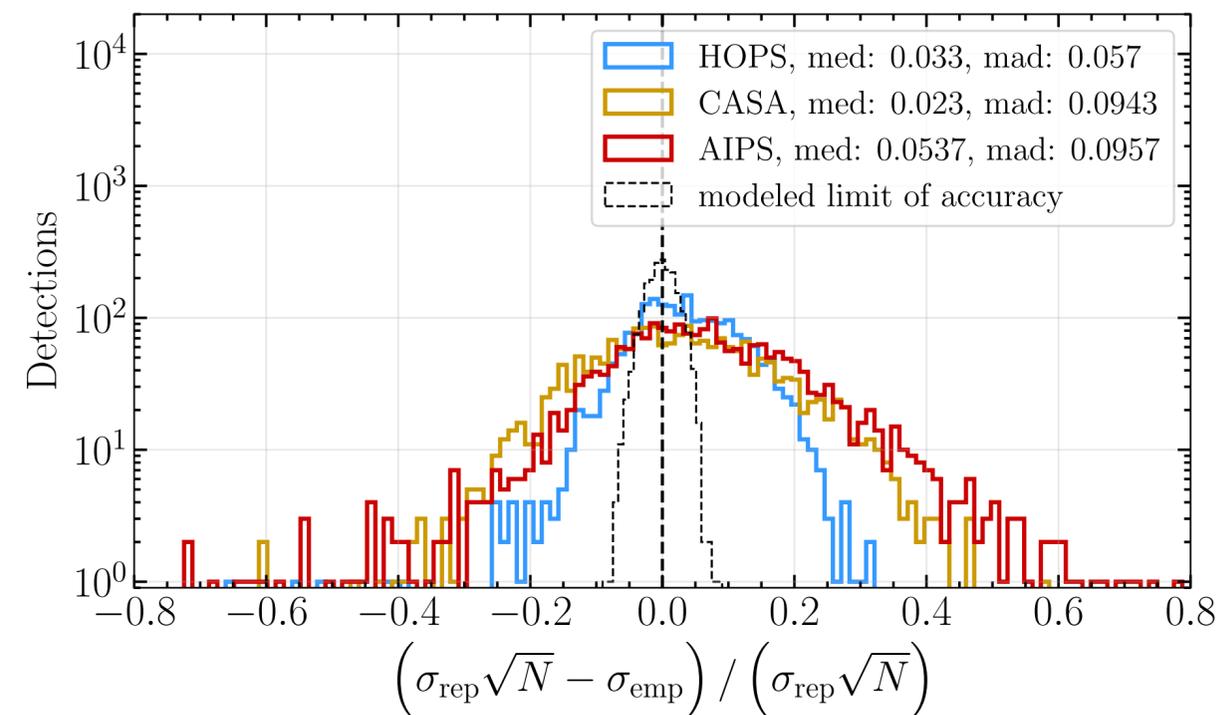
Thermal (statistical) error due to contribution from independent system noise at each site. For a normalized correlation coefficient and white noise, this follows from the **central limit theorem**,

$$\sigma_{r,ij}^2 = \frac{1}{2 \Delta t \Delta \nu}$$

Thermal noise is **Gaussian** and **independent** in real, imaginary components, and thus scales very simply under vector average and scaling by any visibility amplitude factors. Still, it is always good to check!



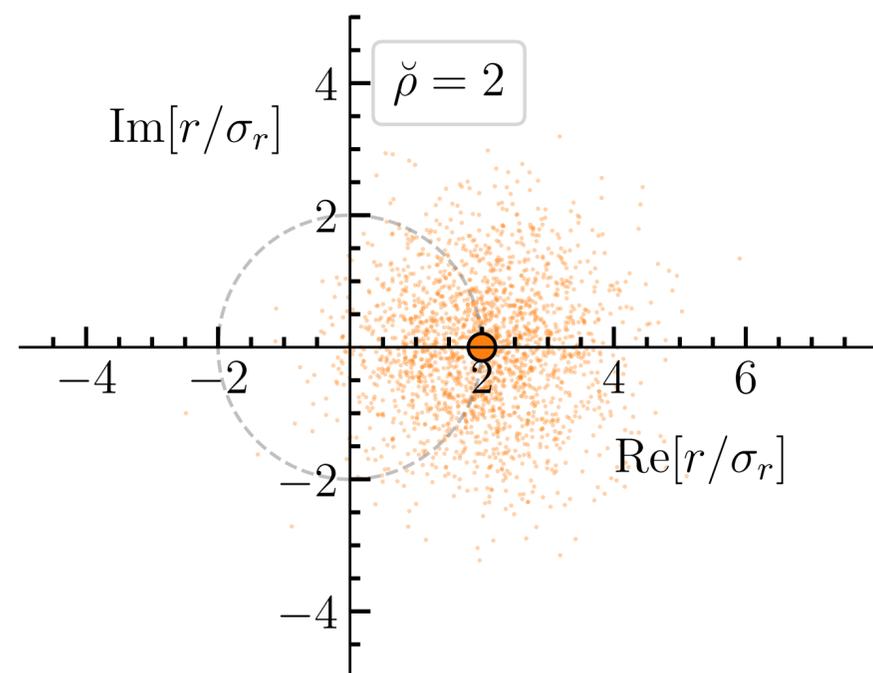
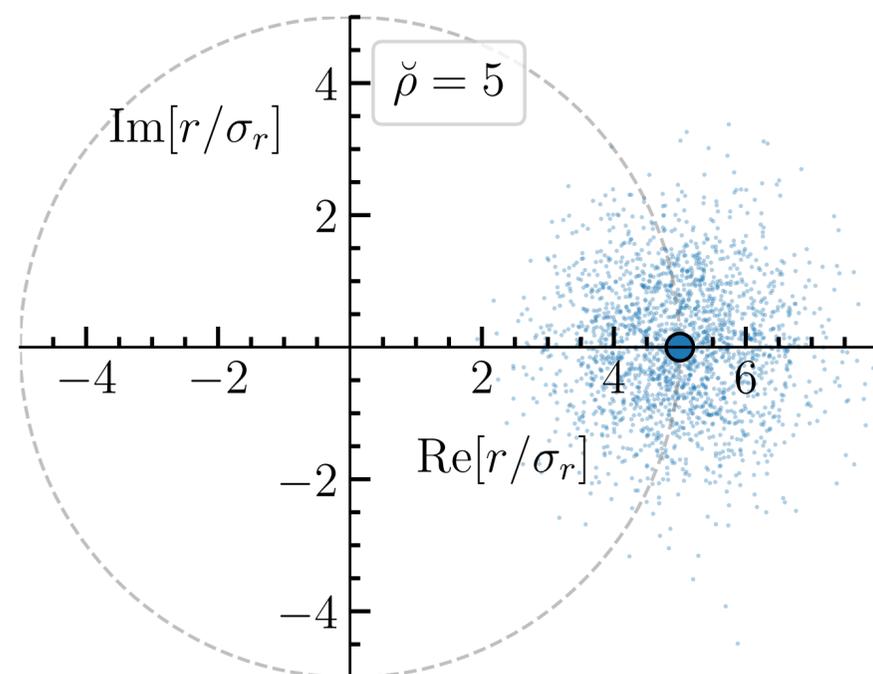
“Closure-phase” differencing, e.g. Ortiz+ 2016



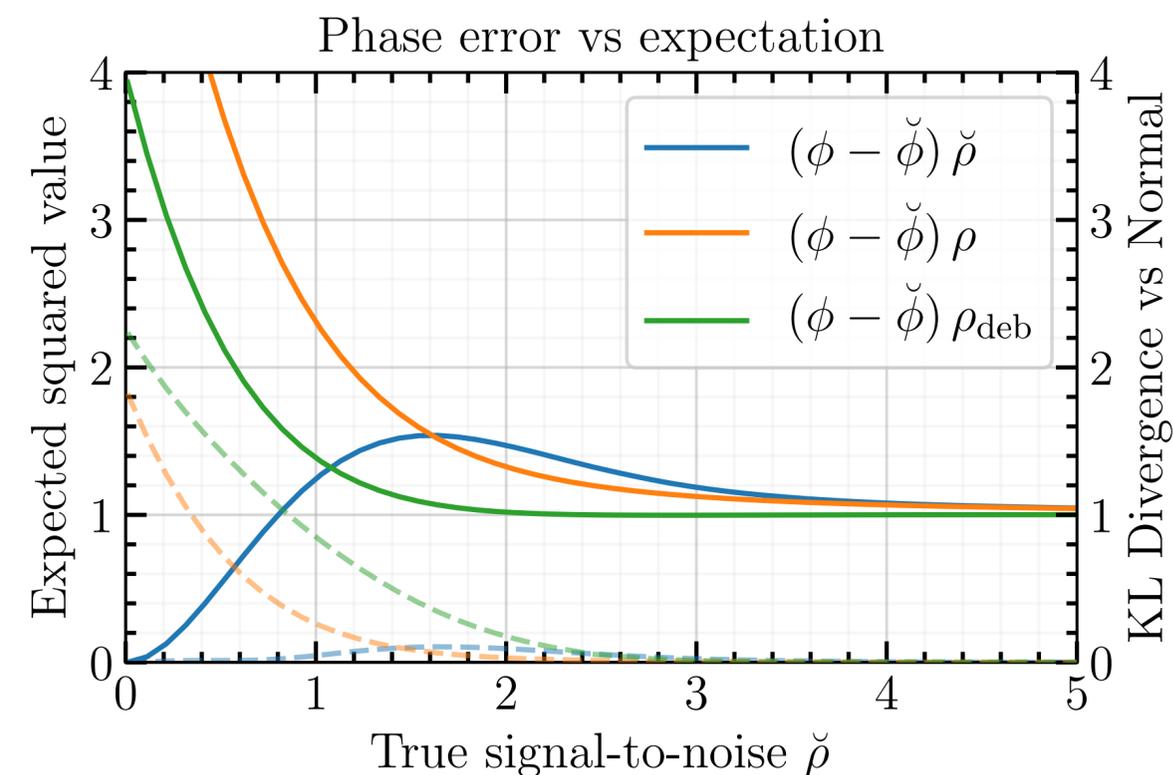
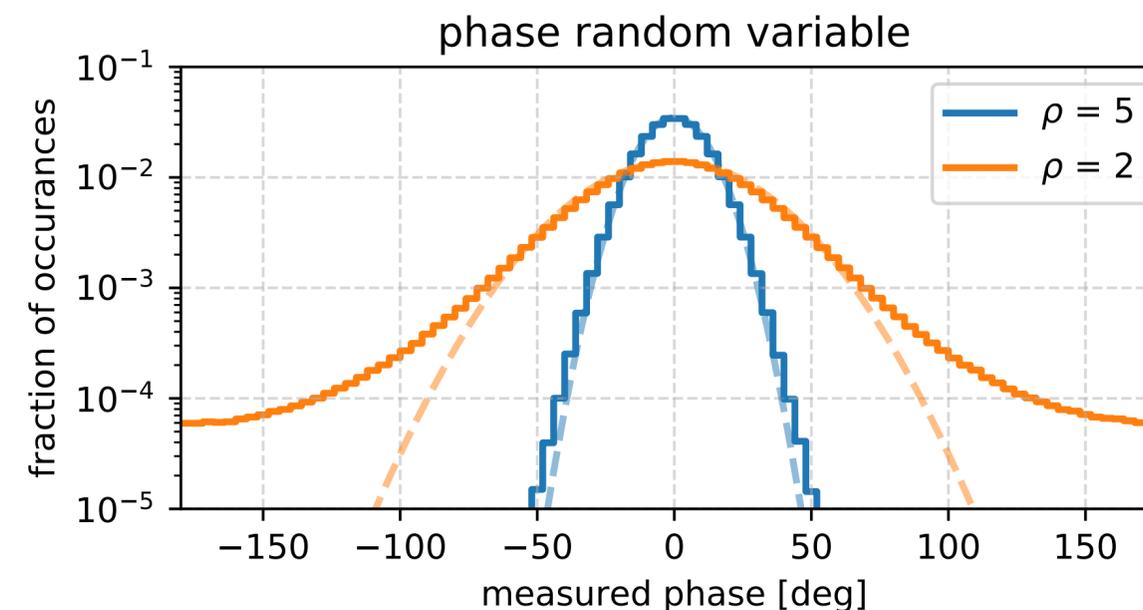
Amplitude scatter, e.g. Wielgus+ 2019-CE-02

Thermal errors: non-Gaussianity

Thermal error is Gaussian in complex visibility, not necessarily in amplitude & phase



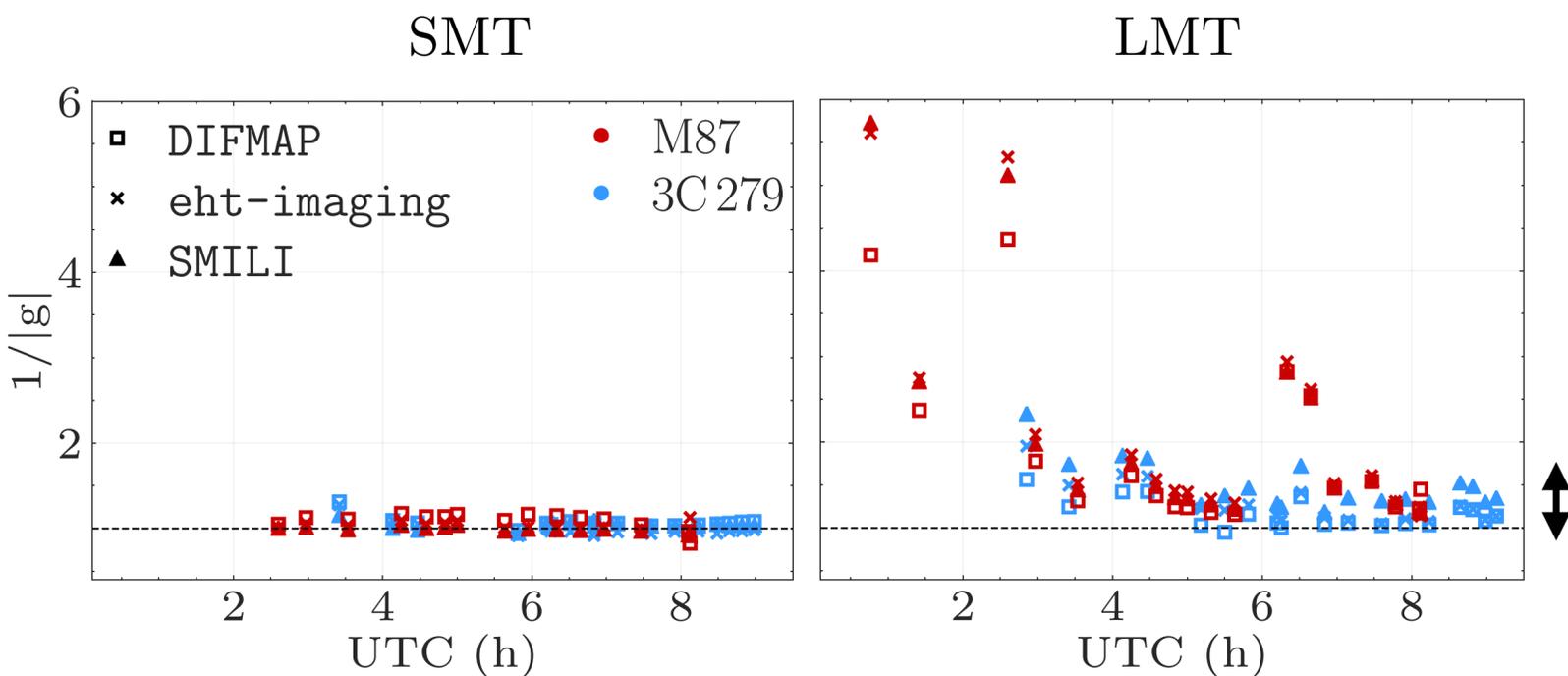
common estimators of phase error will give large reduced chi-square at low S/N



Systematic errors: closing vs non-closing

Closing errors (manageable)

Errors in gain calibration: $V_{ij} = g_i g_j^* r_{ij} + n_{ij}$

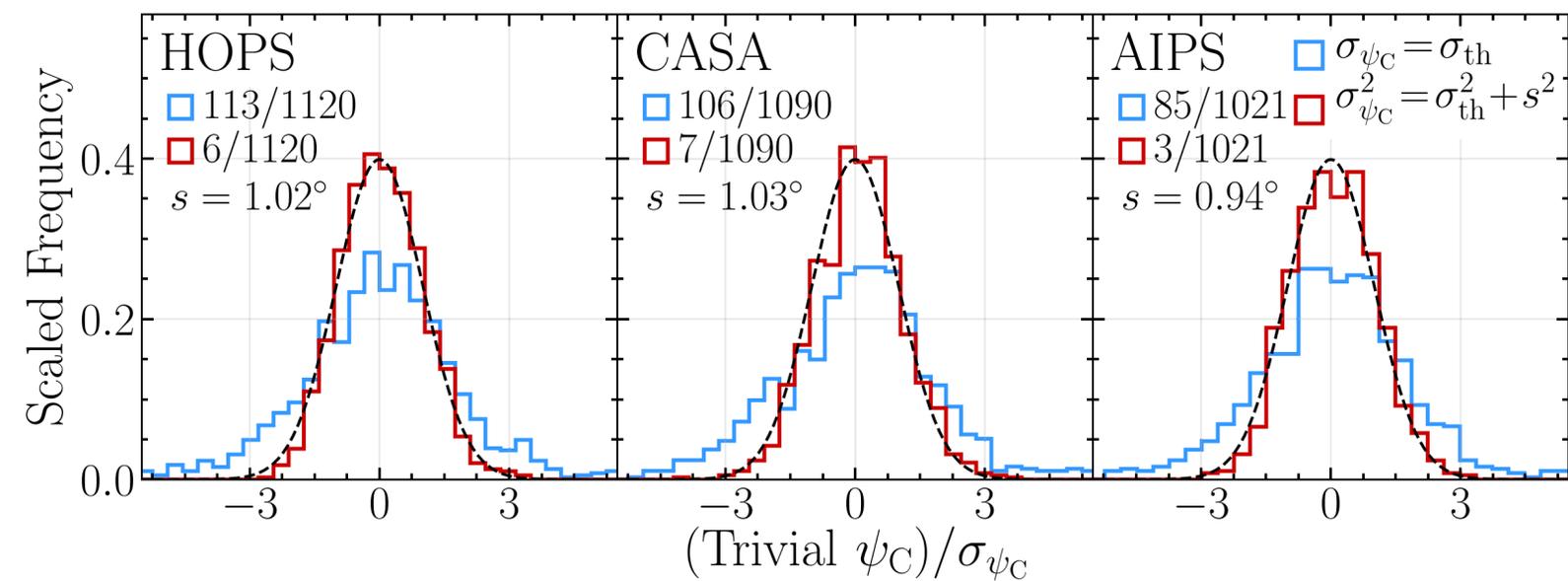


Possibly reflected in high/low band comparison, pipeline comparison, etc

If uncertain, best left for self-calibration (do not “inflate” data errors)

Non-closing errors (try to minimize)

Non-thermal baseline errors: $V_{ij} = g_i g_j^* r_{ij} + n_{ij} + e_{ij}$



Difficult to estimate, possibly reflected in trivial closure phases and amplitudes polarization leakage, band-pass non-overlap, coherence issue, etc

Commonly modeled as additional Gaussian RV:

$$\sigma^2 = \sigma_{\text{th}}^2 + s^2 \quad s \sim 1\text{-}2\%$$

but be careful! most likely not independent across data points (do not average..)



Covariant errors

The noise properties of the correlation coefficients from the correlator are very simple:

Gaussian noise in real and imaginary components, **independent** across all data products

This is **ideal** for model fitting, calculating likelihoods, goodness-of-fit, etc.. messing with the data just makes it worse

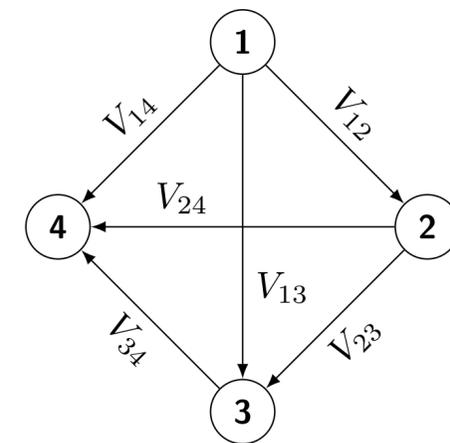
Simple example — Gain error: $V_{ij} = \boxed{g_i} g_j^* r_{ij} + n_{ij}$

If Gaussian, can be captured by covariance matrix (e.g. for log amplitude)

$$\Sigma_a = \begin{pmatrix} \sigma_{12}^2 + \sigma_{g,1}^2 + \sigma_{g,2}^2 & \sigma_{g,1}^2 & \sigma_{g,2}^2 \\ \sigma_{g,1}^2 & \sigma_{13}^2 + \sigma_{g,1}^2 + \sigma_{g,3}^2 & \sigma_{g,3}^2 \\ \sigma_{g,2}^2 & \sigma_{g,3}^2 & \sigma_{23}^2 + \sigma_{g,2}^2 + \sigma_{g,3}^2 \end{pmatrix}$$

More complicated example — Closure phase

$$\Sigma_\psi = \begin{pmatrix} \sigma_{12}^2 + \sigma_{23}^2 + \sigma_{13}^2 & \sigma_{12}^2 & -\sigma_{13}^2 \\ \sigma_{12}^2 & \sigma_{12}^2 + \sigma_{24}^2 + \sigma_{14}^2 & \sigma_{14}^2 \\ -\sigma_{13}^2 & \sigma_{14}^2 & \sigma_{13}^2 + \sigma_{34}^2 + \sigma_{14}^2 \end{pmatrix}$$



covariance over 3 closure phases

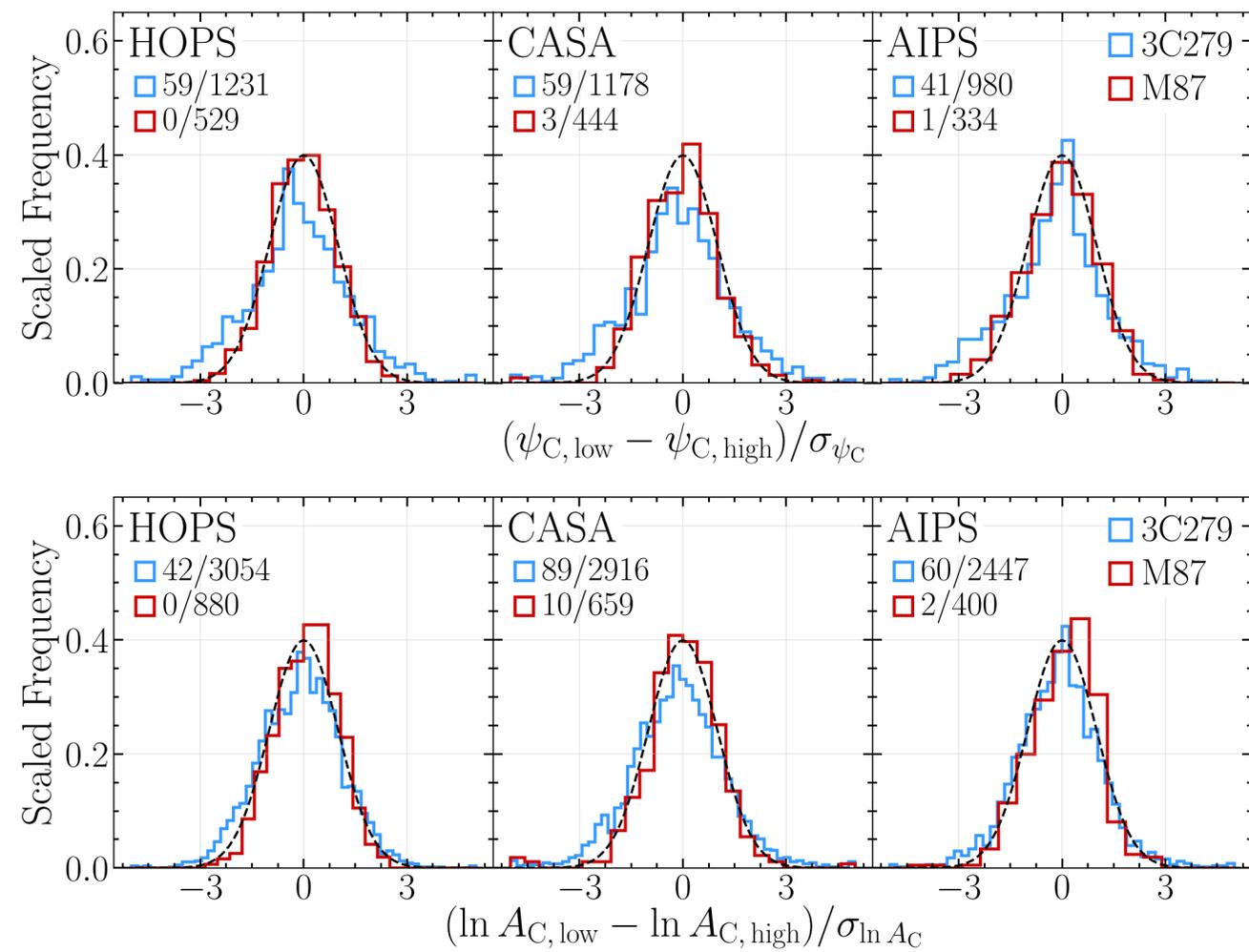
Blackburn, Pesce+ 2019

Ignoring covariant errors often leads to confidence intervals which are too small!

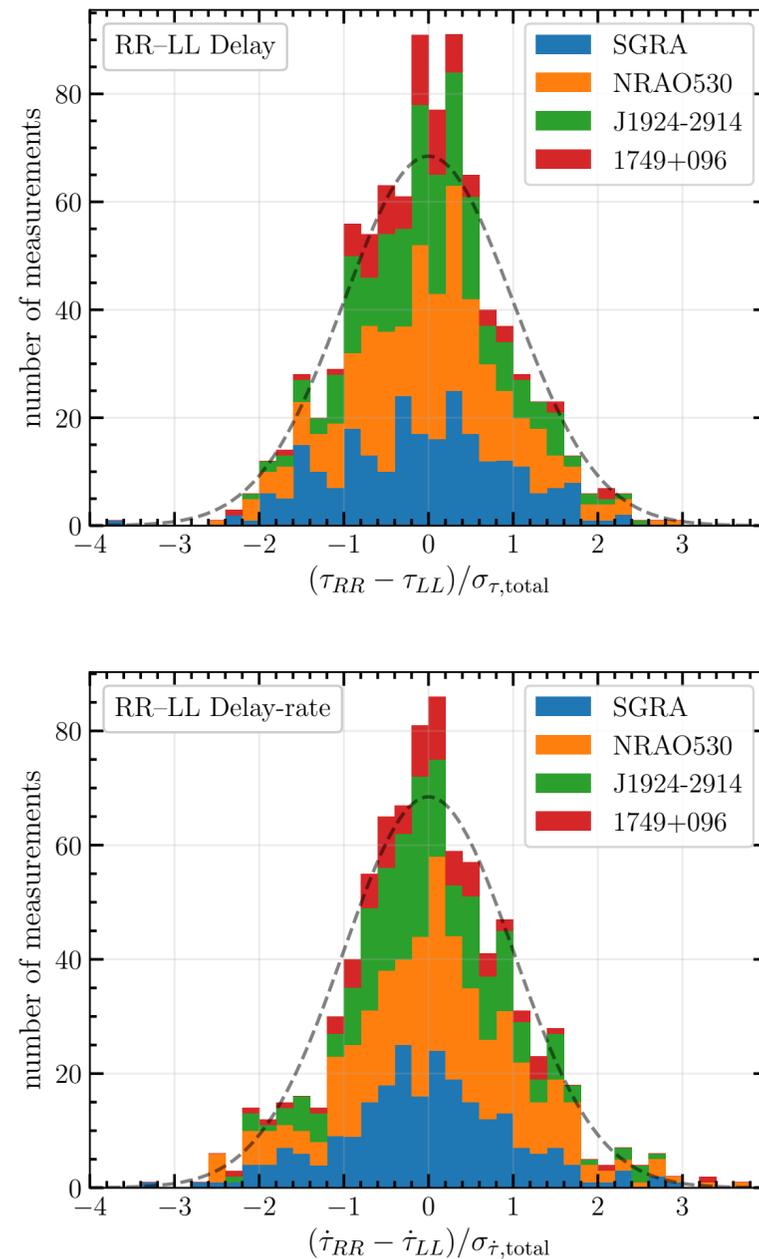
It can be fun and instructive to use covariant errors, but make sure there is a **very good reason** before moving away from forward modeling into **simple data products**..

After a lot of work, we want to make sure we have good data, and not bad data!

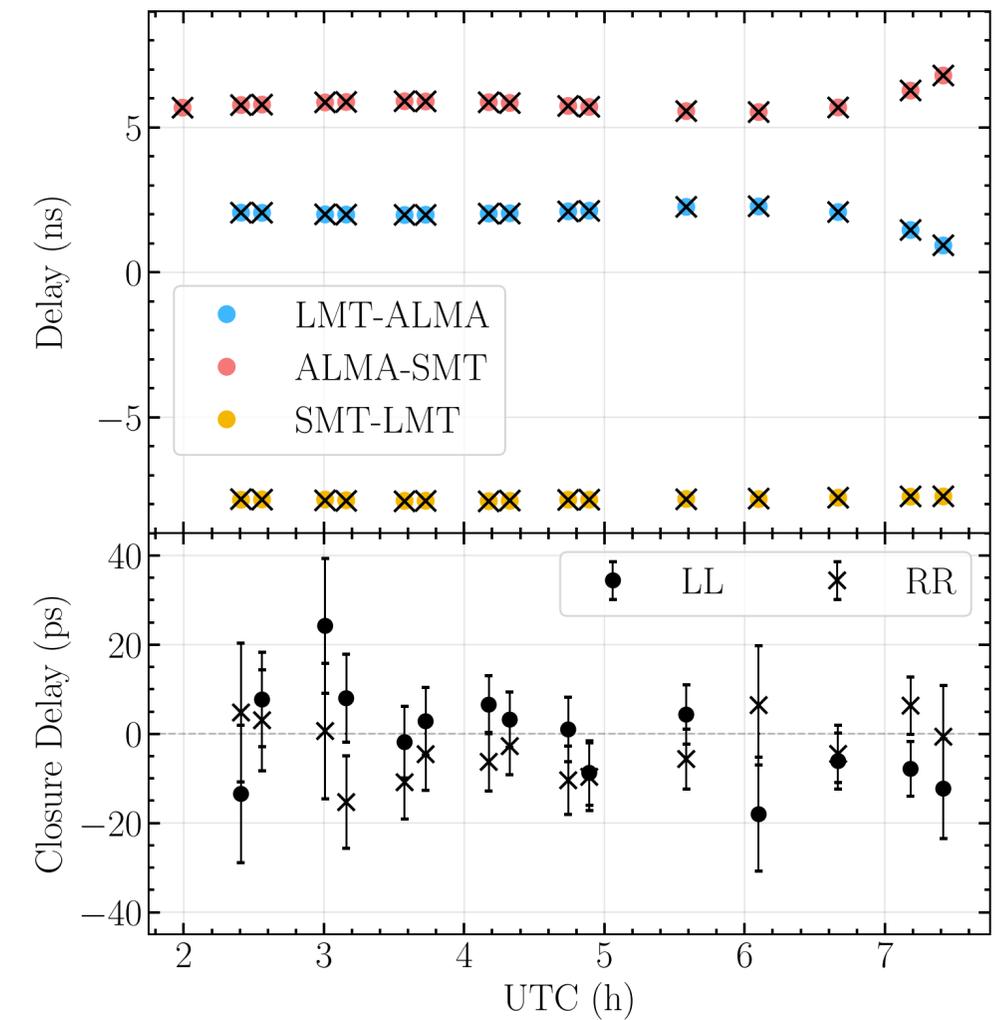
Closure quantity & pipeline cross-comparison



Self-consistency of baseline-based delay solutions



Baseline delay closure





<http://bit.ly/HandlingDataEval>