

# Comparison of Linear and Circular Products after PolConvert

L. V. Benkevitch<sup>1</sup>

<sup>1</sup>MIT Haystack observatory, Westford, MA 01886, USA.

August 30, 2024

## Abstract

The VGOS data are collected with the use of linearly polarized receivers. However, circular polarization has a number of benefits. The PolConvert software allows VGOS data conversion from linear to circular polarization. In this study we do a statistical comparison of the pseudo-Stokes data before and after the conversion in order to characterize the errors PolConvert introduces.

## Contents

<b>1</b>	<b>Parameter Temporal Variations Before and After PolConvert</b>	<b>1</b>
<b>2</b>	<b>Parameter residual Statistics Before and After PolConvert</b>	<b>3</b>
<b>3</b>	<b>Conclusion</b>	<b>4</b>
<b>4</b>	<b>Appendix: Software</b>	<b>4</b>
4.1	make_sorted_idx.py: Saving VGOS data in Python dictionaries . . . . .	4
<b>5</b>	<b>Appendix: Figures</b>	<b>6</b>

## 1 Parameter Temporal Variations Before and After PolConvert

Converting linear polarization of the source data to the circular polarization with the PolConvert software introduces changes into its properties. Here we consider the changes in multi-band delays (MBD), single-band delays (SBD), and signal-to-noise ratios (SNR) for the pseudo-Stokes polarization products, of which we only select the pseudo-Stokes  $I$  parameters as functions of time,  $t$ . From the original, linearly polarized data set we read  $I$ , which we denote  $\text{LinI}(t)$ . PolConvert is applied to the original data set and generates another data set with circular polarization, from which we read  $I$  parameters denoted as  $\text{CirI}(t)$ . Of course, the data are presented in the discrete, numerical format, so instead of the continuous time,  $t$ , integer indices are used, like  $k = \overline{1..N}$ ,

where  $N$  is the number of data points in the sample. The following Figures show temporal variations of the  $\text{LinI}[k]$  and  $\text{CirI}[k]$  parameters during the experiment before and after PolConvert, one graph for each baseline:

MBD variations in Fig. 1;  
SBD variations in Fig. 4;  
SNR variations in Fig. 7.

One can see that most of the graphs after PolConvert,  $\text{CirI}[k]$ , are significantly biased with respect to  $\text{LinI}[k]$ . The biases calculated as  $\text{LinI}[k] - \text{CirI}[k]$  are plotted in the Figures below:

MBD bias variations in Fig. 2;  
SBD bias variations in Fig. 5;  
SNR bias variations in Fig. 8.

However, removing the bias by subtracting the mean values of each curve makes the differences much less significant:

$$\text{LinI}[k] = \text{LinI}[k] - \text{mean}(\text{LinI}[k]) \quad (1)$$

$$\text{CirI}[k] = \text{CirI}[k] - \text{mean}(\text{CirI}[k]) \quad (2)$$

In order to compare the unbiased  $\text{LinI}[k]$  with  $\text{CirI}[k]$ , the residual,  $\text{res}[k]$ , is computed as their difference:

$$\text{res}[k] = \text{LinI}[k] - \text{CirI}[k] \quad (3)$$

The graphs of residuals with the means subtracted are shown in the following Figures:

MBD residual variations in Fig. 3;  
SBD residual variations in Fig. 6;  
SNR residual variations in Fig. 9.

To assess the similarity of  $\text{LinI}[k]$  and  $\text{CirI}[k]$ , the following parameters are computed:

- Pearson correlation coefficient  $r_{corr}$ ;
- Root Mean Square Error, RMSE;

Since  $\text{LinI}[k]$  and  $\text{CirI}[k]$  already have their means subtracted, the correlation formula is simplified:

$$r_{corr} = \frac{\sum \text{LinI}[k] \cdot \text{CirI}[k]}{\sqrt{\sum \text{LinI}[k]^2 \cdot \sum \text{CirI}[k]^2}} \quad (4)$$

The root mean square error, RMSE, in our case is actually the standard deviation of the residual,  $\sigma$ :

$$\text{RMSE} = \sqrt{\frac{1}{N} \sum \text{res}[k]^2} \quad (5)$$

Lower RMSE indicates a better fit between the curves. However, RMSE is not an “absolute” indicator of the curves’ proximity. The correlation coefficient, though, has its upper limit, unity, and if  $r_{corr}$  is very close to unity, the curves are almost identical, save their constant bias. MBD and SNR demonstrate excellent correlations between LinI and CirI.

## 2 Parameter residual Statistics Before and After PolConvert

While in the previous Section MBD, SBD, and SNR pseudo-Stokes  $I$  parameter data were treated as functions of time, here we do their statistical study. In order to estimate the errors introduced by PolConvert, histograms of the unbiased parameter residuals (after the mean subtraction, see Eqs. (1) and (2)) are plotted. Below are the histograms for all of the baselines:

MBD residual histogram for all of the baselines in Fig. 10;  
 SBD residual histogram for all of the baselines in Fig. 12;  
 SNR residual histogram for all of the baselines in Fig. 14.

It is interesting to see which of the individual stations contribute to the errors. We plotted histograms of the parameter residuals for the baselines involving each particular station. They are shown in the following Figures:

MBD residual histograms for the baselines including one station in Fig. 11;  
 SBD residual histograms for the baselines including one station in Fig. 13;  
 SNR residual histograms for the baselines including one station in Fig. 15.

Note that everywhere the means of the residuals are very close to zero (actually, down to the machine epsilon!).

In order to evaluate the significance of the error we attempted to use the Pearson’s chi-squared test comparing the histograms to the normal distributions:

$$f(\mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}. \quad (6)$$

All the histograms of residuals are created with 21 bins. To perform the Pearson’s  $\chi^2$  test, the left-tail and right-tail bins with sparse data (frequencies less than 5) were grouped and the test used fewer number of bins  $N_b$  (see it printed in each plot, Fig. 10 - Fig. 15). The Pearson’s  $\chi^2$  test compares the observed data counts in the bins  $B_i$  with the expected (theoretical) counts  $E_i$ , where  $i = \overline{1 \dots N_b}$ , obtained from the normal PDF (6) as  $E_i = Np_i$ . Here  $p_i$  is the probability that a random value sampled from  $f(\mu, \sigma)$  is within the  $i$ -th interval. The value of the test-statistic is calculated as

$$X^2 = \sum_{i=1}^{N_b} \frac{(B_i - E_i)^2}{E_i}, \quad (7)$$

where  $N_b$  is the number of bins. The  $X^2$  test statistic asymptotically approaches the  $\chi^2_{df}$  distribution, where df is the number of degrees of freedom:

$$df = N_b - k - 1, \quad (8)$$

where  $k = 2$  is the number of parameters estimated from the data, mean  $\mu$  and standard deviation  $\sigma$ . The additional  $-1$  is due to the fact that the  $\chi^2$  test itself has a restriction imposed by the sum of the observed frequencies being equal to the sum of the expected frequencies.

If we want to be 95% confident that the data is distributed normally,  $X^2$  cannot exceed the critical value  $\chi_{cr}^2$  such that the probability for  $X^2$  to appear within the interval  $[0.. \chi_{cr}^2]$  is  $p = 0.95$ . The critical  $\chi_{df}^2$  value can be calculated as the value of  $\chi^2$  Probability Point Function or PPF, the inverse of the  $\chi^2$  Cumulative Distribution Function or CDF. In Python this is done as `scipy.stats.chi2.ppf(1-alpha, df)`, where  $\alpha = 0.05$  is the level of significance.

However, in our case the Pearson's  $\chi^2$  test does not work, because the residuals are grouped around their means substantially denser than the normal distributions with the same standard deviations. In most distributions, the calculated values of the histogram chi-squares are many times greater than the critical chi-square values for the p-values less than or equal to 0.05 (printed in each of the histogram plots).

We have used another method. Since the residuals are mostly within  $\pm\sigma$  (standard deviation), we compute their proportion and compare it to the expected proportion under a normal distribution, which is about 68.27%. The proportion is printed near each histogram. It is at the level of about 80%, which is better than the standard normal for most of the cases.

### 3 Conclusion

We have conducted a study to assess how well the PolConvert software makes VGOS data conversion from linear to circular polarization. The pseudo-Stokes data of the Multi-Band Delays (MBD), Single-Band Delays (SBD), and Signal-to-Noise Ratios (SNR) as functions of time before and after the conversion were compared. We also did a statistical comparison of the pseudo-Stokes data of MBD, SBD, and SNR before and after the conversion in order to characterize the errors PolConvert introduces. Two criteria were used to assess the residuals: the Pearson's  $\chi^2$  test and the proportions of residuals within  $\pm\sigma$  compared to the normal distribution.

PolConvert introduces large (and random!) biases between the pseudo-Stokes data parameters, which, we believe, is unacceptable.

On the other hand side, after the bias removal, the average errors (computed as unbiased residual rms) are very small making tenths of percent.

## 4 Appendix: Software

### 4.1 `make_sorted_idx.py`: Saving VGOS data in Python dictionaries

The VLBI Global Observing System (VGOS) database is organized as a tree-like directory structure. For our purpose of statistical analysis of a small number of parameters scattered across many directories and files below the root directory of the experiment, this implies significant overhead in opening multiple files and accessing the parameters within each of them. The data files in their names only provide the station or baseline names, and no time or polarization information. For example, extraction of, say, SNR data for a particular polarization product and within a specific time range would require opening *all* the files and accessing their times and polarizations using HOPS API calls.

We wrote a script, `make_sorted_idx.py`, to extract the parameters for statistical analysis for the whole experiment and to put it in a Python dictionary, preserving the temporal order. We

call such dictionaries “indices”. The index can be “pickled” and saved on disk. Interestingly, these files are small, in the hundreds of kilobytes. The other data analysis and plotting scripts read the index files, unpickle them into the Python dictionaries, and use data from the dictionaries.

The script `make_sorted_idx.py` should be run on the `demi.haystack.mit.edu` server where the VGOS data are stored under the directory `/data-sc16/geodesy/`. Current version of the script works on the experiment 3819 data located under `/data-sc16/geodesy/3819`. It creates three dictionaries pickled in the files `idx3819l.pkl`, `idx3819c.pkl`, and `idx3819cI.pkl` in the directory where the script was run.

- `idx3819l.pkl`: linear polarization products immediately from the `/data-sc16/geodesy/3819/` directory;
- `idx3819c.pkl`: circular polarization products generated by PolConvert without the pseudo-Stokes 'I' data, only 'LL', 'LR', 'RL', 'RR'.  
The data are found in `/data-sc16/geodesy/3819/polconvert/3819/scratch/pol_prods1/3819` directory.
- `idx3819cI.pkl`: pseudo-Stokes 'I' only for the circular polarization products generated by PolConvert. The data are taken from `/data-sc16/geodesy/3819/polconvert/3819/scratch/pcphase_stokes_test/3819`

A pickle file can be unpickled into a dictionary using the `pickle.load()` function. For example:

```
import pickle
with open('idx3819c.pkl', 'rb') as finp:
    idx3819c_1 = pickle.load(finp)
```

The script `make_sorted_idx.py` is also a Python module defining the function

`make_idx(base_dir, pol='lin', max_depth=2)` with parameters:

`base_dir`: the directory containing the VGOS data. For example, it may be `/data-sc16/geodesy/3819/`.

`pol`: polarization, 'lin' - linear, 'cir' - circular.

This parameter is used for the data generated by PolConvert.

It converts the polarization product names

'XX', 'XY', 'YX', 'YY', and the lists ['XX', 'YY'] into the correct names 'LL', 'LR', 'RL', 'RR', and 'I', respectively.

`max_depth`: Limits the maximum depth of recursing into the subdirectories of `base_dir`.

`make_idx()` creates and returns the index dictionary with the data from `base_dir`.

The index dictionary has three dimensions: the baseline name, the polarization, and the data proper, including 'time', 'file', 'mbdelay', 'sbdelay', and 'snr'. Consider a particular index named `idx3819l_1` (experiment 3819, linear polarization). Its first dimension is indexed with the baseline names derived from the set of stations, {'E', 'M', 'S', 'T', 'V', 'Y'}.

The possible first indices are the baseline names:

```
idx38191_1.keys()
dict_keys(['SE', 'VY', 'MV', 'MT', 'TV', 'EY', 'SY', 'TY', 'MS', 'SV',
          'TE', 'EV', 'MY', 'ME'])
```

Each of the baselines is associated with the cross-correlation products and the pseudo-Stokes I parameter. Thus the second index is one of the products. For example, for the 'ME' baseline:

```
idx38191_1['ME'].keys()
dict_keys(['XX', 'XY', 'YX', 'YY', 'I'])
```

For example, the times, the full data file names, SNRs, multi- and single-band delays for the 'SY' baseline and the 'XY' polarization products from this baseline are contained in the index dictionary under

```
idx38191_1['SY']['XY']:
```

```
idx38191_1['SY']['XY'].keys() prints
dict_keys(['time', 'file', 'mbdelay', 'sbdelay', 'snr']).
```

For example, in order to access the multi-band delay data list in the ascending temporal order for the baseline 'SV' and the pseudo-Stokes I, one should issue the following command:

```
mbd = idx38191_1['SV']['I']['mbdelay'].
```

## 5 Appendix: Figures

Pseudo-Stokes  $I$  MBD (ps) vs Time (min), Lin & Cir Pol after PolConvert

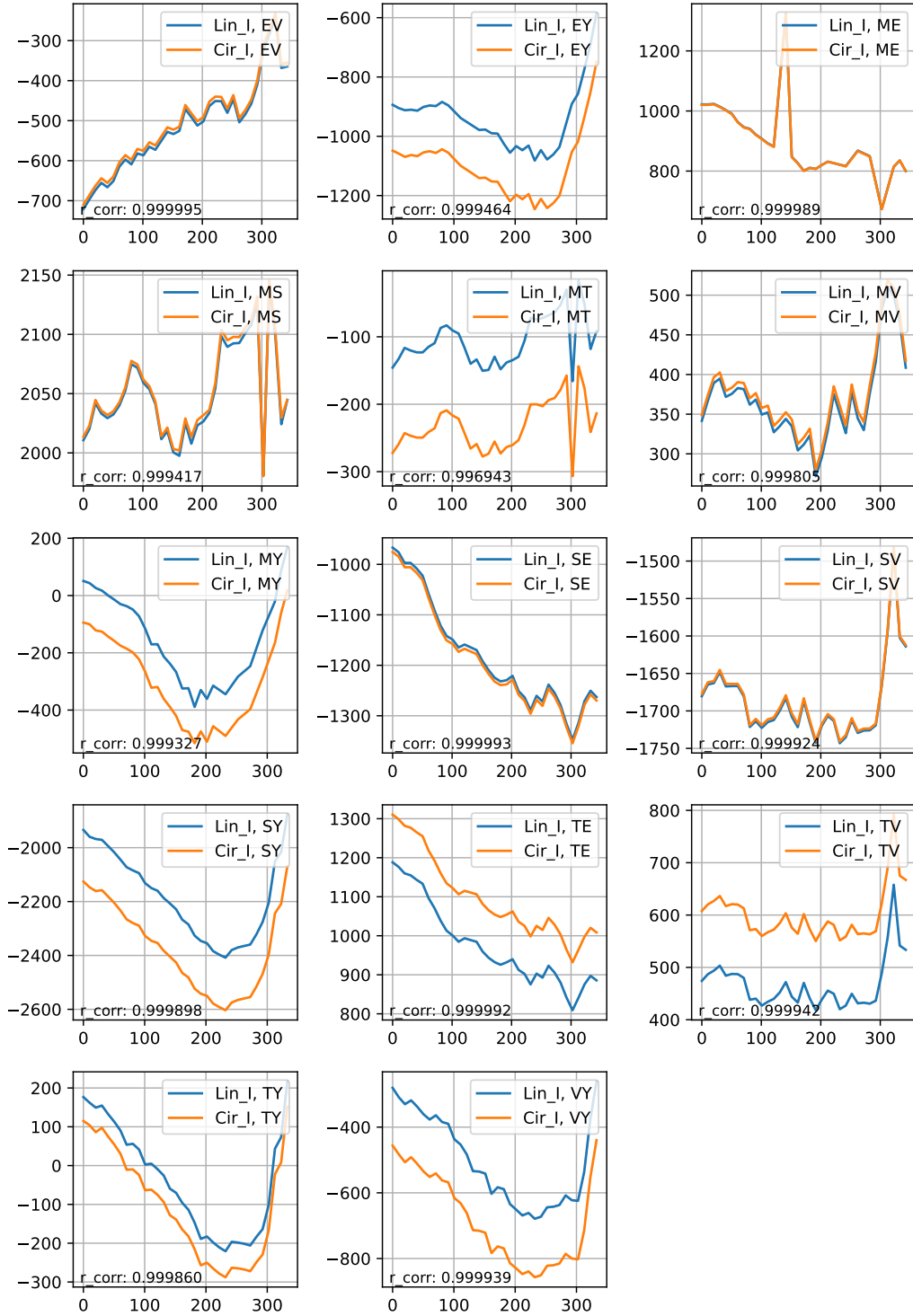


Figure 1: Evolution of Multi-Band Delays (in picoseconds) during the experiment (time in minutes) for the pseudo-Stokes  $I$  parameter for every baseline. Each panel shows two curves: LinI,  $I$  from the original, linearly polarized data set and CirI,  $I$  from the circularly polarized data set obtained by applying PolConvert to the original data. Graphs have large biases. However, the coefficients of correlation  $r\_corr$  between LinI and CirI are so close to unity that the curves can be considered identical (except for the biases).

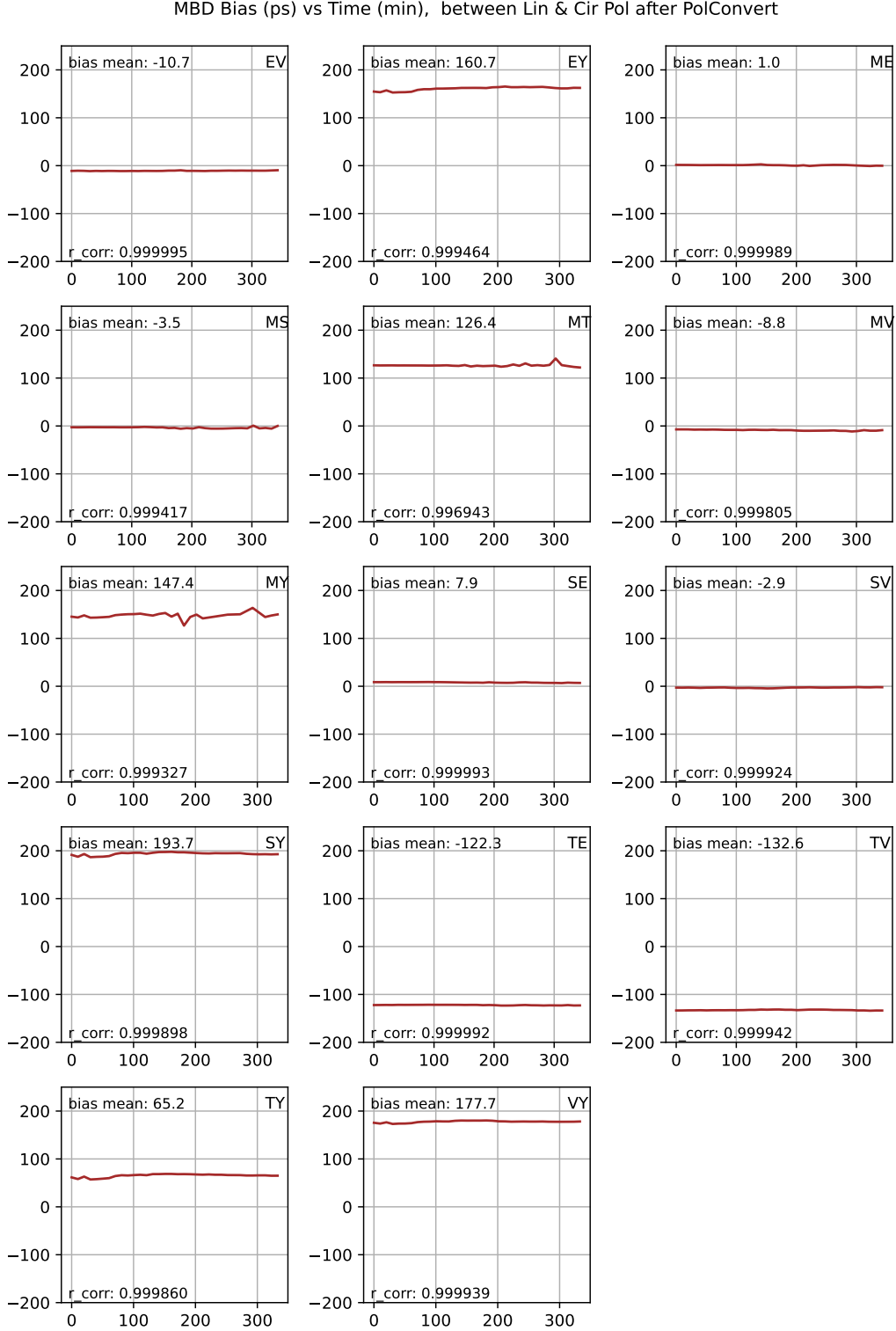


Figure 2: Evolution of the bias between LinI and CirI Multi-Band Delays (in picoseconds) during the experiment (time in minutes) for the pseudo-Stokes  $I$  parameter for every baseline. Each panel shows one curve of LinI - CirI, where LinI is  $I$  from the original, linearly polarized data set and CirI is  $I$  from the circularly polarized data set obtained by applying PolConvert to the original data. The means for each curve and the coefficients of correlation  $r\_corr$  between LinI and CirI are given.



MBD Residuals (ps) vs Time (min), between Lin & Cir Pol after PolConvert (means subtracted)

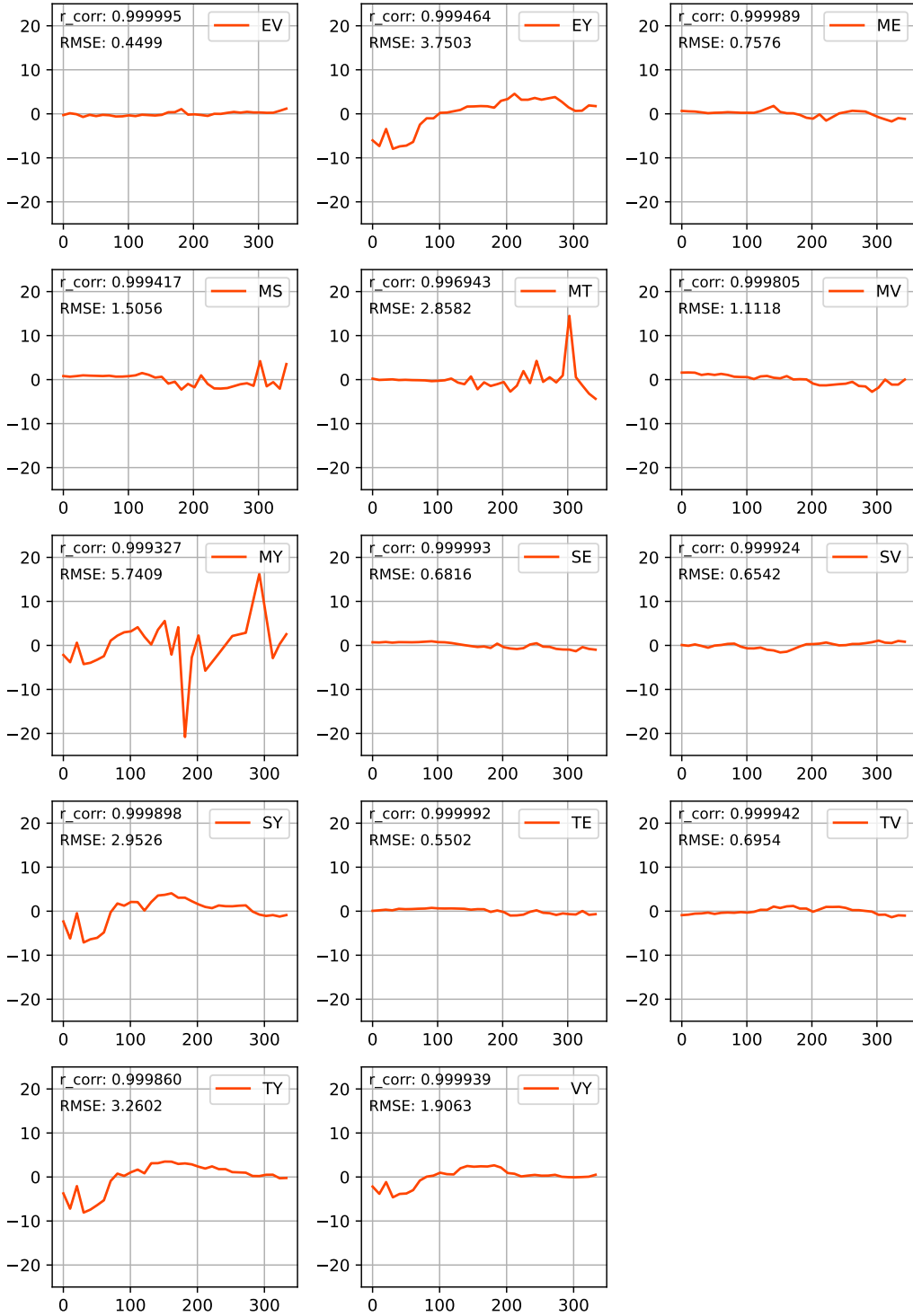


Figure 3: Evolution of the residuals between unbiased LinI and CirI Multi-Band Delays (in picoseconds) during the experiment (time in minutes) for the pseudo-Stokes  $I$  parameter for every baseline. Each panel shows a curve of LinI - CirI, both with subtracted means, where LinI is  $I$  from the original, linearly polarized data set and CirI is  $I$  from the circularly polarized data set obtained by applying PolConvert to the original data. The coefficients of correlation  $r\_corr$  between LinI and CirI are very close to unity. The Root Mean Square Error (RMSE) shows the standard deviation of the residuals.

Pseudo-Stokes  $I$  SBD (ps) vs Time (min), Lin & Cir Pol after PolConvert

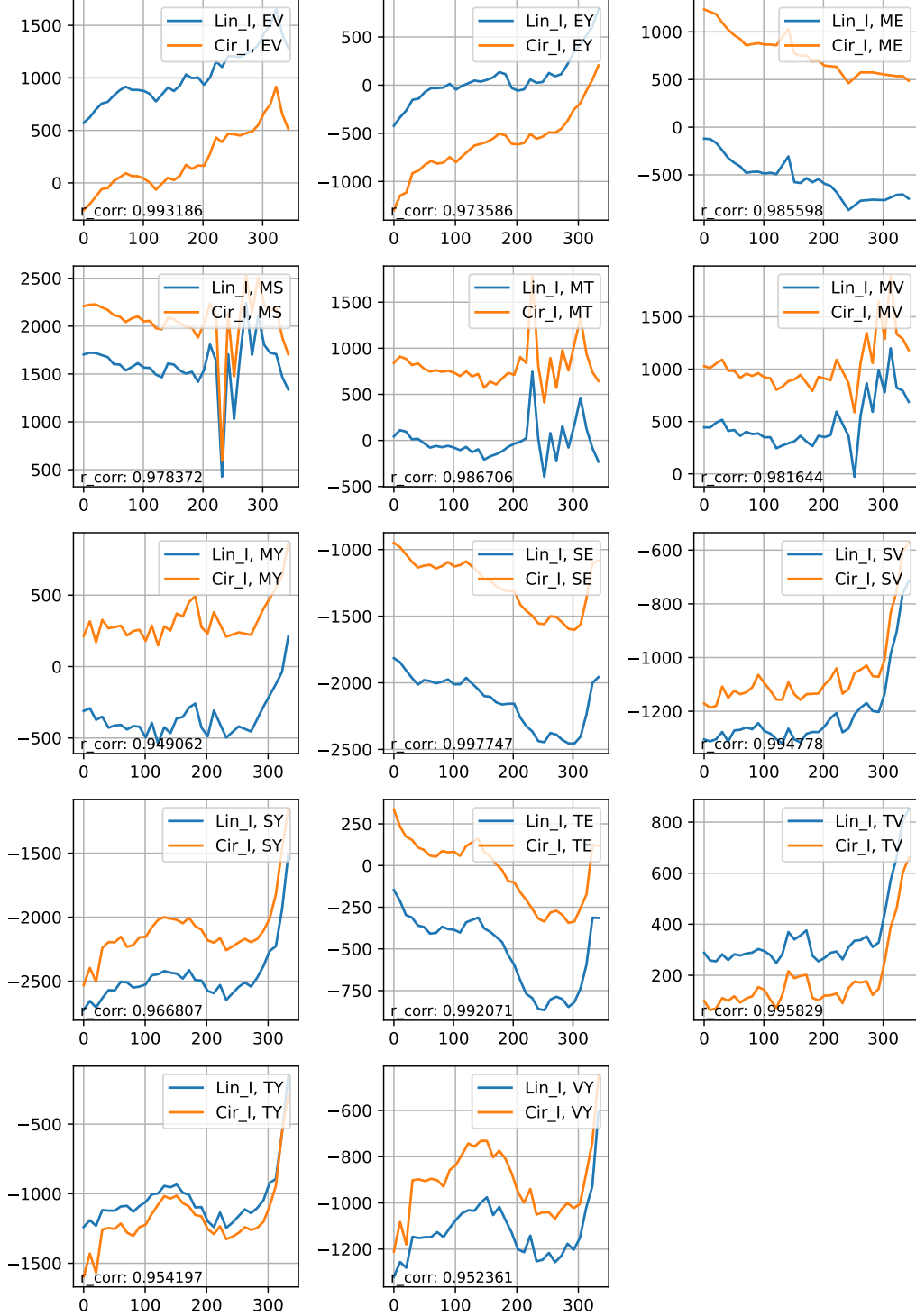


Figure 4: Evolution of Single-Band Delays (inpicoseconds) during the experiment (time in minutes) for the pseudo-Stokes  $I$  parameter for every baseline. Each panel shows two curves: LinI,  $I$  from the original, linearly polarized data set and CirI,  $I$  from the circularly polarized data set obtained by applying PolConvert to the original data. Graphs have large biases. The coefficients of correlation  $r\_corr$  between Lin\_I and CirI for most of the baselines are close to unity. If it were not for the biases, the curves could be considered close to identical.

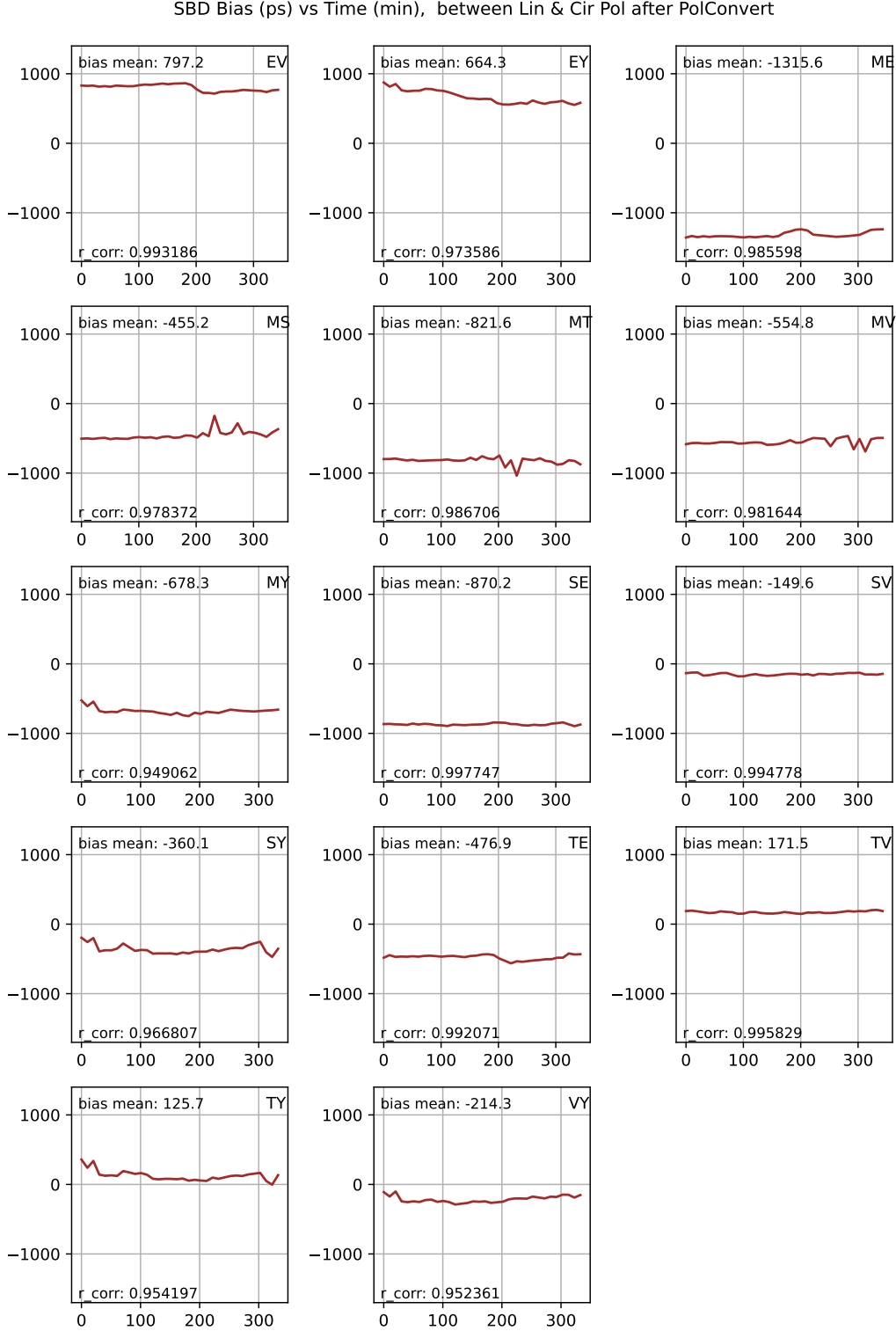


Figure 5: Evolution of the bias between LinI and CirI Single-Band Delays (in picoseconds) during the experiment (time in minutes) for the pseudo-Stokes  $I$  parameter for every baseline. Each panel shows a curve of LinI - CirI, where LinI is  $I$  from the original, linearly polarized data set and CirI is  $I$  from the circularly polarized data set obtained by applying PolConvert to the original data. The means for each curve and the coefficients of correlation  $r\_corr$  between LinI and CirI are given.

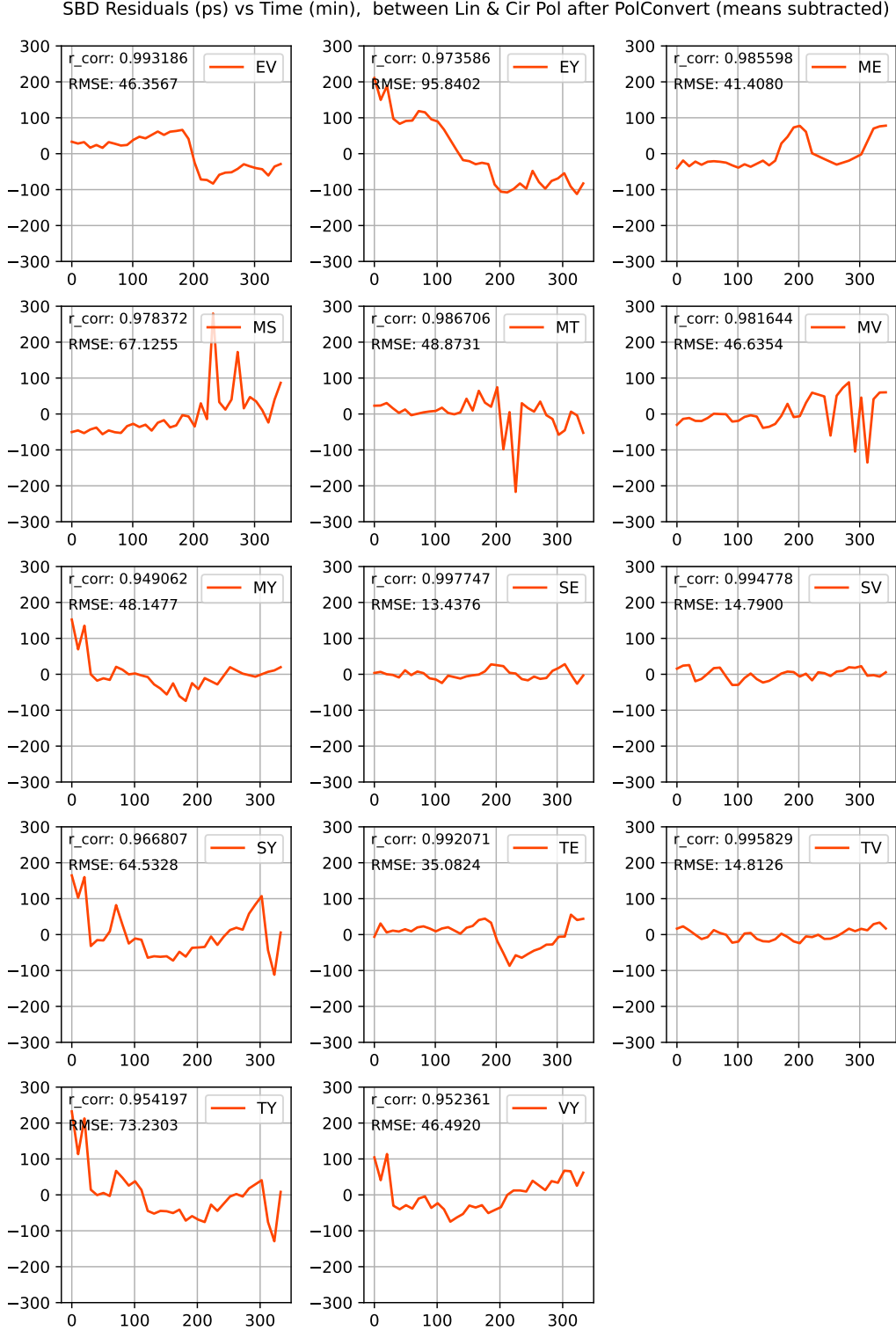


Figure 6: Evolution of the residuals between unbiased LinI and CirI Single-Band Delays (in picoseconds) during the experiment (time in minutes) for the pseudo-Stokes  $I$  parameter for every baseline. Each panel shows a curve of LinI - CirI, both with subtracted means, where LinI is  $I$  from the original, linearly polarized data set and CirI is  $I$  from the circularly polarized data set obtained by applying PolConvert to the original data. The coefficients of correlation  $r\_corr$  between LinI and CirI are reasonably close to unity. The Root Mean Square Error (RMSE) shows the standard deviation of the residuals.

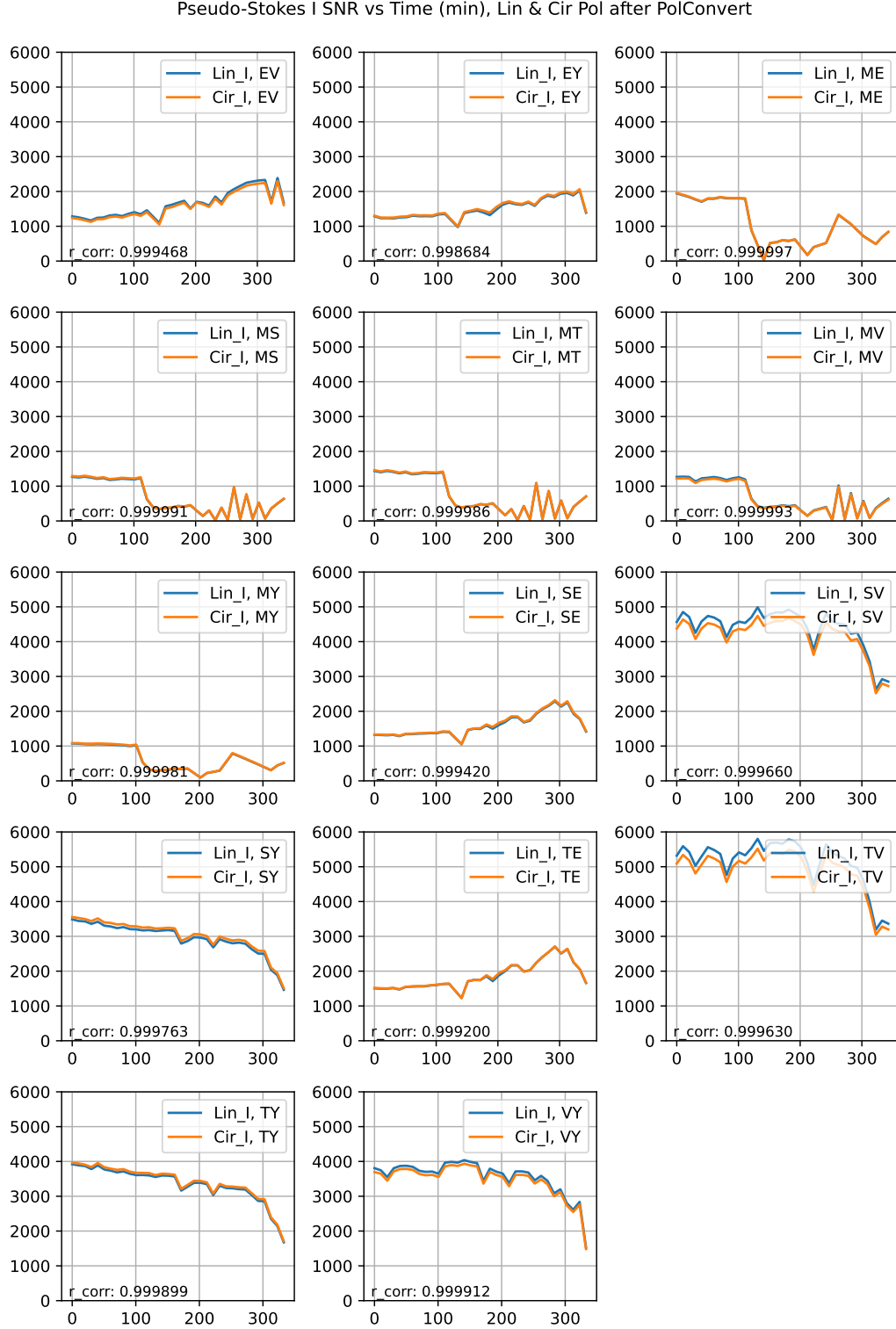


Figure 7: Evolution of Signal-to-Noise Ratios during the experiment (time in minutes) for the pseudo-Stokes  $I$  parameter for every baseline. Each panel shows two curves: LinI,  $I$  from the original, linearly polarized data set and CirI,  $I$  from the circularly polarized data set obtained by applying PolConvert to the original data. Graphs have large biases. However, the coefficients of correlation  $r_{\text{corr}}$  between LinI and CirI are so close to unity that the curves can be considered identical (except for the biases).

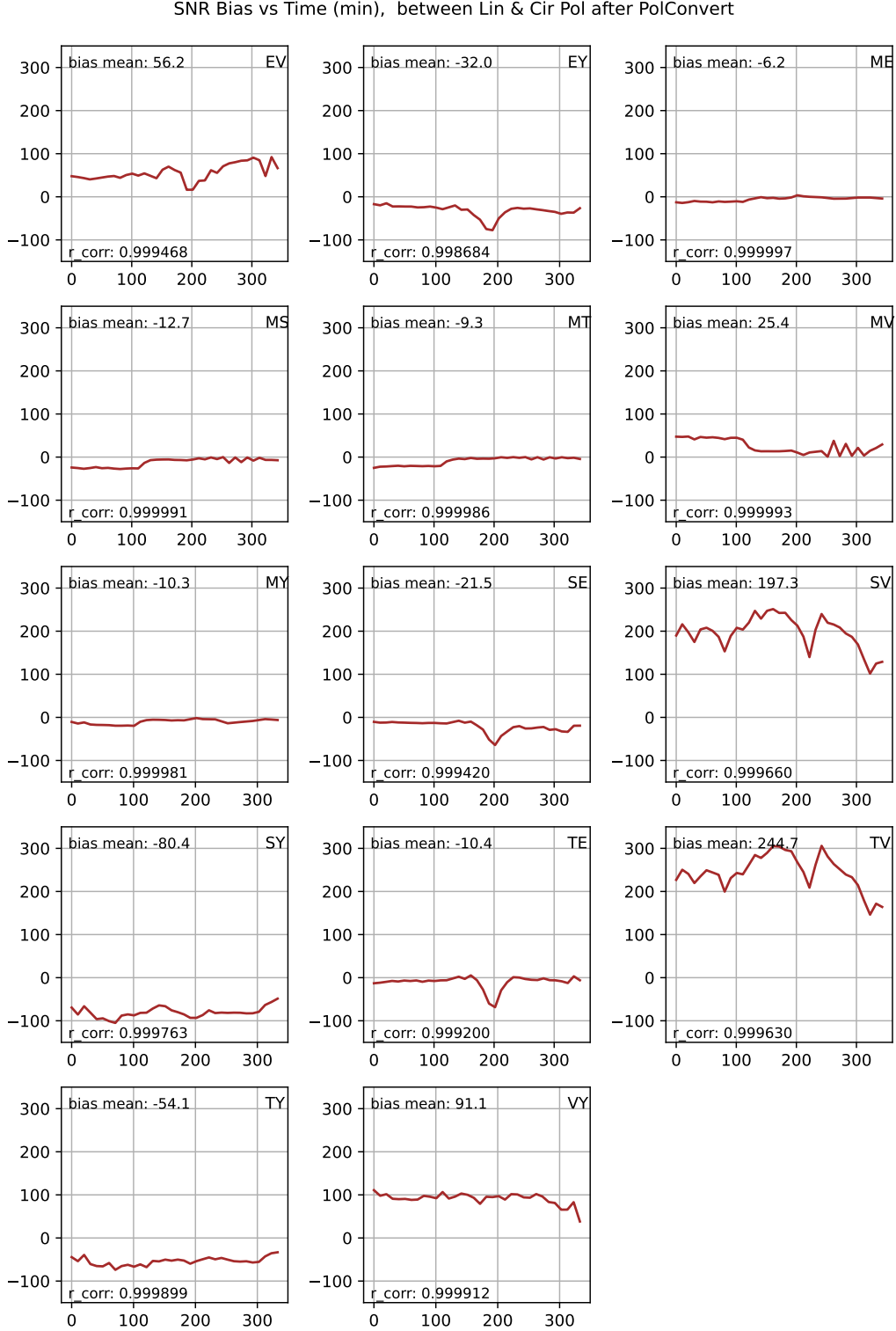


Figure 8: Evolution of the bias between LinI and CirI Signal-to-Noise Ratios during the experiment (time in minutes) for the pseudo-Stokes  $I$  parameter for every baseline. Each panel shows a curve of LinI - CirI, where LinI is  $I$  from the original, linearly polarized data set and CirI is  $I$  from the circularly polarized data set obtained by applying PolConvert to the original data. The means for each curve and the coefficients of correlation  $r\_corr$  between LinI and CirI are given.

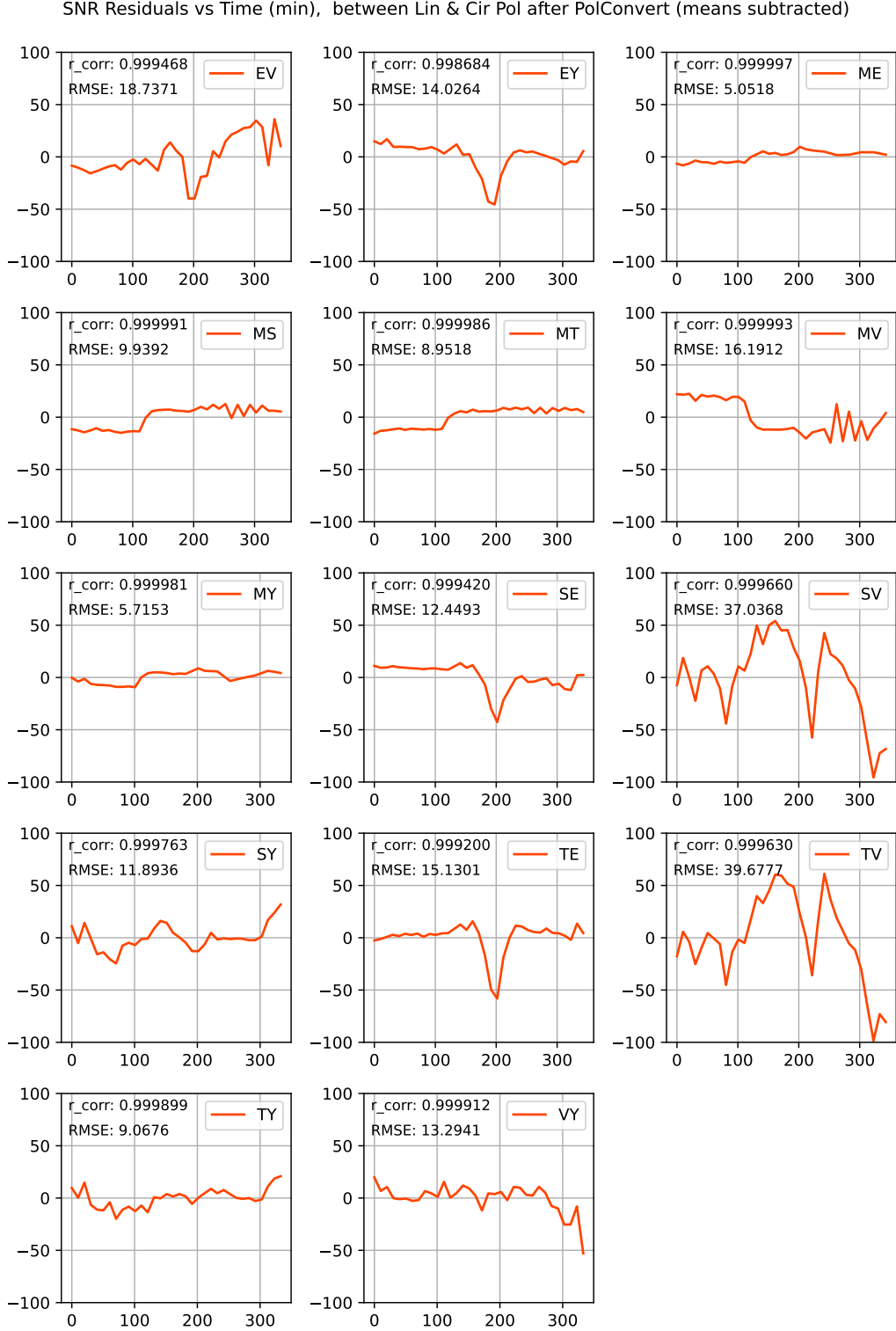


Figure 9: Evolution of the residuals between unbiased LinI and CirI Signal-to-Noise Ratios during the experiment (time in minutes) for the pseudo-Stokes  $I$  parameter for every baseline. Each panel shows a curve of LinI - CirI, both with subtracted means, where LinI is  $I$  from the original, linearly polarized data set and CirI is  $I$  from the circularly polarized data set obtained by applying PolConvert to the original data. The coefficients of correlation  $r_{\text{corr}}$  between LinI and CirI are very close to unity. The Root Mean Square Error (RMSE) shows the standard deviation of the residuals.

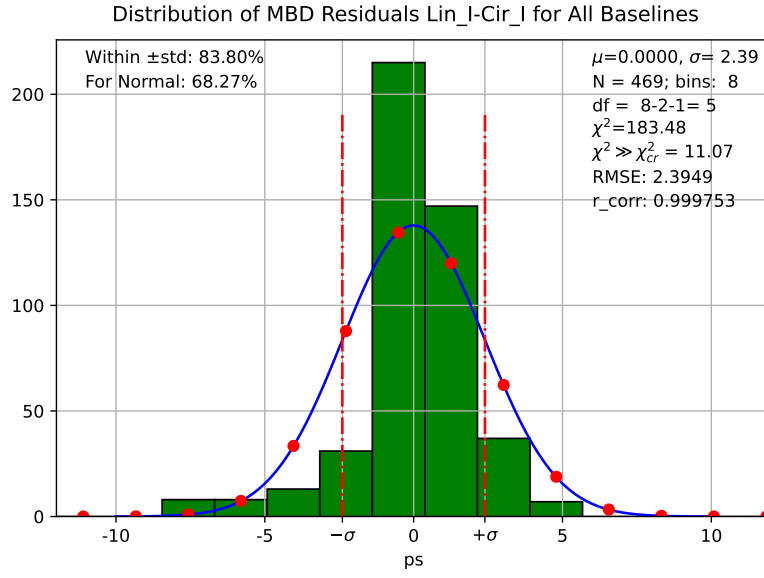


Figure 10: Histogram of the Multi-Band Delay residuals (see Eq. (3)) for all of the baselines.  $N$  is the total number of residuals spread over 21 bin. For the Pearson's  $\chi^2$  test, the left-tail and right-tail bins with sparse data (less than 5 counts in each) were grouped to make only  $N_b = 8$  bins. The number of degrees of freedom  $\text{df}$  is 5 (i.e.  $N_b - 3$ ). With this  $\text{df}$ , the critical value,  $\chi^2_{cr}$  is 11.07. The  $\chi^2$  observed, 183.48, greatly exceeds  $\chi^2_{cr}$ , so the test fails. However, the proportion of residuals within  $\pm\sigma$ , 83.8%, is much better than 68.27% of a normal distribution.



Distributions of MBD Residuals Lin\_I-Cir\_I for Stations

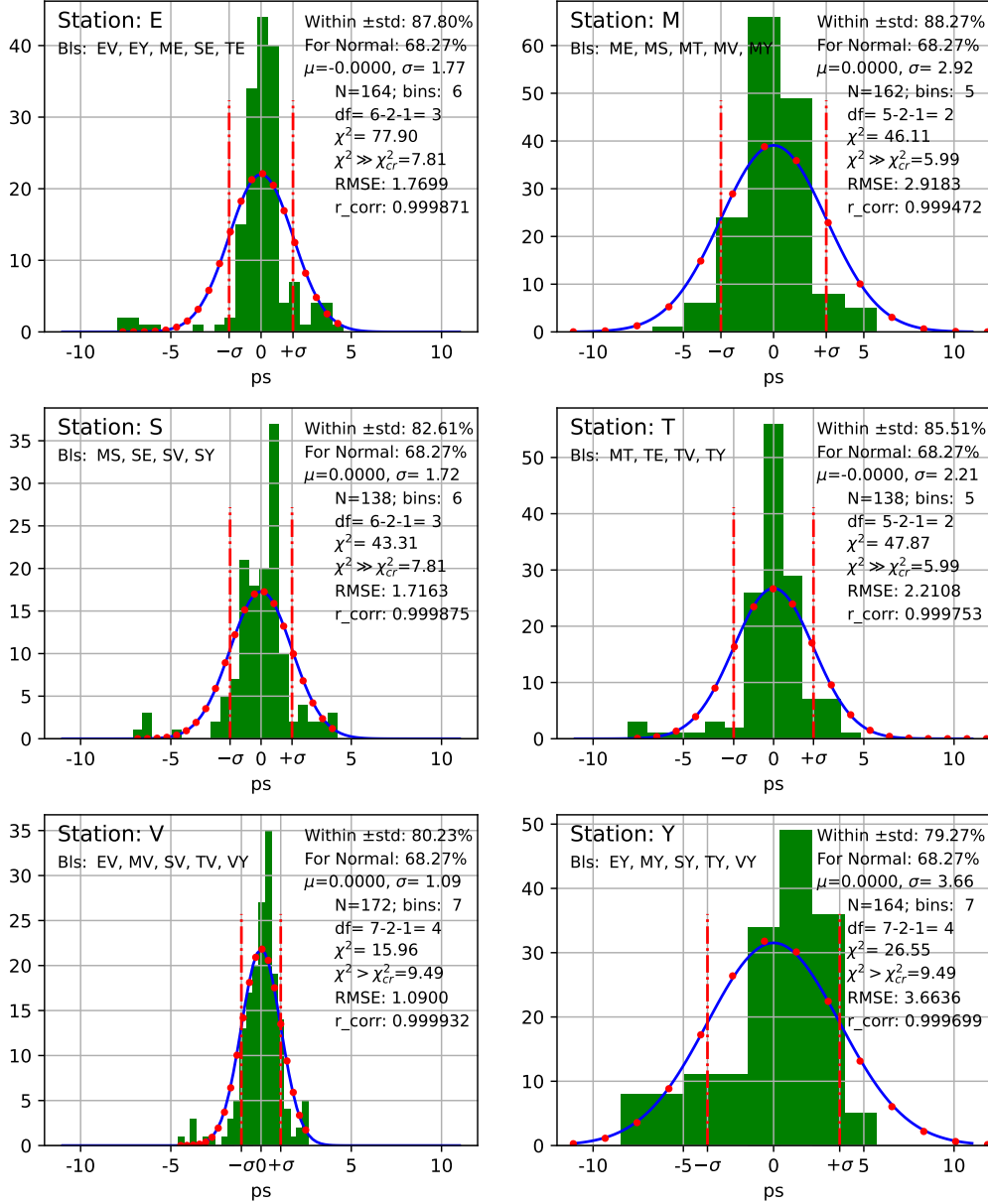


Figure 11: Histograms of the Multi-Band Delay residuals (see Eq. (3)) for the groups of baselines including an individual station. N is the total number of residuals spread over 21 bin. For the Pearson's  $\chi^2$  test, the left-tail and right-tail bins with sparse data (less than 5 counts in each) were grouped to make smaller number of bins  $N_b$ . The number of degrees of freedom df is  $N_b - 3$ . With these df, the critical values  $\chi^2_{cr}$  are calculated. For all the distributions, the  $\chi^2$  observed exceed  $\chi^2_{cr}$ , so the tests fail. However, the proportions of residuals within  $\pm\sigma$  are much better than 68.27% of a normal distribution.

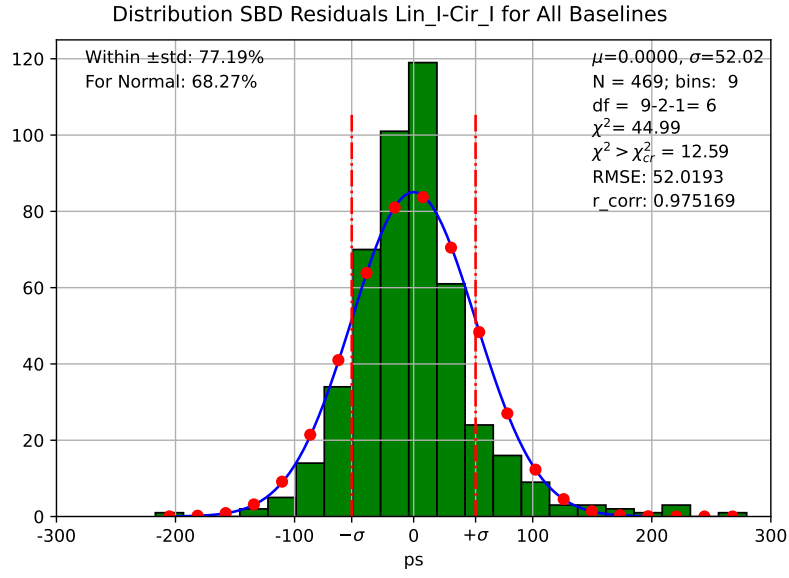


Figure 12: Histogram of the Single-Band Delay residuals (see Eq. (3)) for all of the baselines.  $N$  is the total number of residuals spread over 21 bin. For the Pearson's  $\chi^2$  test, the left-tail and right-tail bins with sparse data (less than 5 counts in each) were grouped to make only  $N_b = 9$  bins. The number of degrees of freedom  $df$  is 6 (i.e.  $N_b - 3$ ). With this  $df$ , the critical value,  $\chi^2_{cr}$  is 12.59. The  $\chi^2$  observed, 44.99, significantly exceeds  $\chi^2_{cr}$ , so the test fails. However, the proportion of residuals within  $\pm\sigma$ , 77.19%, is much better than 68.27% of a normal distribution.

Distributions of SBD Residuals Lin\_I-Cir\_I for Stations

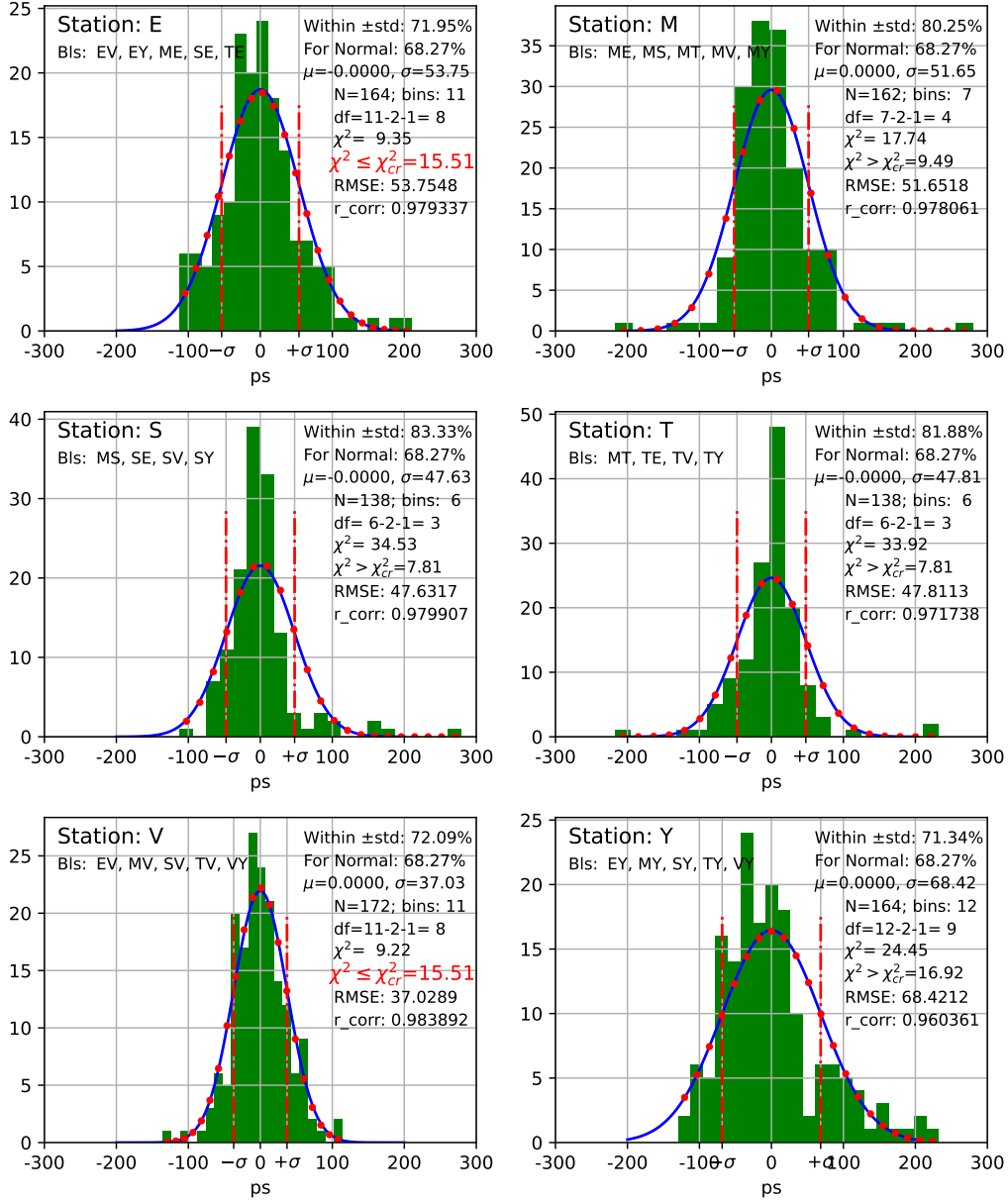


Figure 13: Histograms of the Single-Band Delay residuals (see Eq. (3)) for the groups of baselines including an individual station.  $N$  is the total number of residuals spread over 21 bin. For the Pearson's  $\chi^2$  test, the left-tail and right-tail bins with sparse data (less than 5 counts in each) were grouped to make smaller number of bins  $N_b$ . The number of degrees of freedom  $df$  is  $N_b - 3$ . With these  $df$ , the critical values  $\chi^2_{cr}$  are calculated. For most of the distributions (except those including stations E and V), the  $\chi^2$  observed exceed  $\chi^2_{cr}$ , so the tests fail. Stations E and V have normally distributed residuals with 95% confidence. However, for all of the stations, the proportions of residuals within  $\pm\sigma$  are better than 68.27% of a normal distribution.

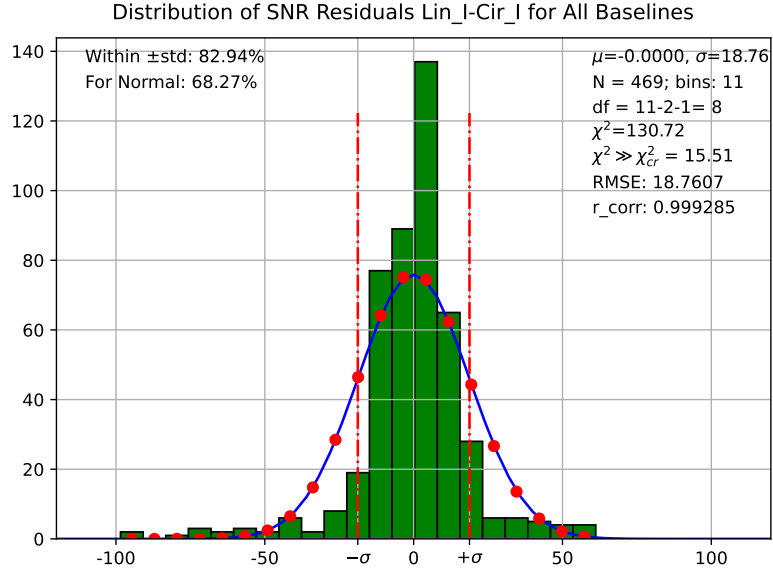


Figure 14: Histogram of the Signal-to-Noise Ratio residuals (see Eq. (3)) for all of the baselines.  $N$  is the total number of residuals spread over 21 bin. For the Pearson's  $\chi^2$  test, the left-tail and right-tail bins with sparse data (less than 5 counts in each) were grouped to make only  $N_b = 11$  bins. The number of degrees of freedom  $\text{df}$  is 8 (i.e.  $N_b - 3$ ). With this  $\text{df}$ , the critical value,  $\chi^2_{cr}$  is 15.51. The  $\chi^2$  observed, 130.72, greatly exceeds  $\chi^2_{cr}$ , so the test fails. However, the proportion of residuals within  $\pm\sigma$ , 82.94%, is much better than 68.27% of a normal distribution.

Distributions of SNR Residuals Lin\_I-Cir\_I for Stations

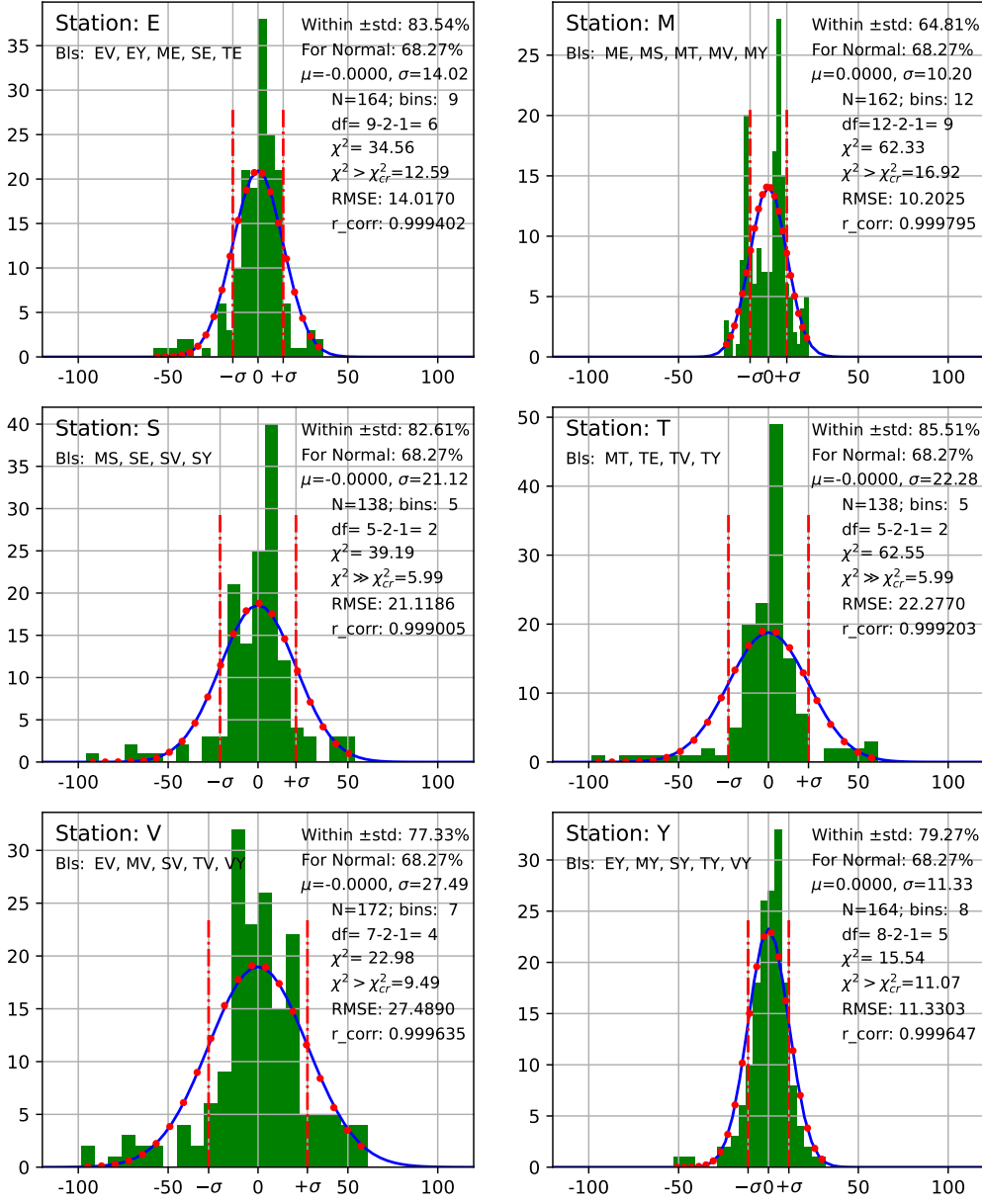


Figure 15: Histograms of the Signal-to-Noise Ratio residuals (see Eq. (3)) for the groups of baselines including an individual station.  $N$  is the total number of residuals spread over 21 bin. For the Pearson's  $\chi^2$  test, the left-tail and right-tail bins with sparse data (less than 5 counts in each) were grouped to make smaller number of bins  $N_b$ . The number of degrees of freedom  $df$  is  $N_b - 3$ . With these  $df$ , the critical values  $\chi^2_{cr}$  are calculated. For all the distributions, the  $\chi^2$  observed exceed  $\chi^2_{cr}$ , so the tests fail. However, the proportions of residuals within  $\pm\sigma$  (except station M with its 64.81%) are much better than 68.27% of a normal distribution.