

Radio Frequency Interference Mitigation Via Cyclostationary Signal Processing: Simulations and Performance Metrics

2 RYAN S. LYNCH 

3 *Green Bank Observatory, PO Box 2, Green Bank, WV, 24944-0002, USA*

4 EVAN T. SMITH 

5 *Green Bank Observatory, PO Box 2, Green Bank, WV, 24944-0002, USA*

6 MATTHEW F. HARRISON

7 *Green Bank Observatory, PO Box 2, Green Bank, WV, 24944-0002, USA*

8 (Received April 22, 2024)

9 Submitted to AJ

10 ABSTRACT

11 We describe an algorithm for identifying radio frequency interference in astronom-
12 ical data by detecting cyclostationarity using the strip spectral correlation analyzer.
13 Cyclostationarity is a property common to many sources of interference but rare in
14 astrophysical sources. We test our algorithm using simulated interfering signals with
15 a variety of modulation processes, symbol durations, numbers of bits-per-symbol, and
16 signal-to-noise ratios, and compare the performance for different algorithmic paramete-
17 rs and thresholds for flagging a signal as interference. We also include a simulated
18 astrophysical spectral line. Our algorithm performs reasonably well for most simulation
19 parameters, with an average area under the resulting receiver operating characteristic
20 curve of 0.90 and ϕ coefficient value of 0.61 when averaged over all signal properties
21 and when using optimal algorithmic parameters. However, we find better performance
22 for subsets of the simulated signals, especially when the signals have relatively narrow
23 bandwidth compared to a spectrometer channel. Our approach does not perform as
24 well for wide-bandwidth signals and frequency-switched signals with large frequency
25 deviations. We discuss potential strategies for improving performance for these types
26 of interferers. We believe cyclostationary signal processing is a promising approach to
27 interference mitigation that can complement other methods.
28
29
30

1. INTRODUCTION

Radio frequency interference (RFI) is a ubiquitous problem in radio astronomy, analogous to light pollution at optical wavelengths. Sources of RFI are legion, including (but certainly not limited to) telecommunications, wireless Internet, navigational aides such as radar and Global Positioning System (GPS), high-speed electronics, and electrical generators and transmission lines. RFI degrades the quality of astronomical data by raising the effective noise floor, sometimes making it impossible to detect weak astrophysical sources, and in extreme cases can damage the sensitive electronics used in modern radio telescopes. Radio astronomy observatories are often built in remote locations, taking advantage of terrain to shield telescopes, and are sometimes protected by regulatory restrictions on the types and strength of nearby transmitters. However, the growing number of satellite transmitters and mobile electronic devices, coupled with ever more sensitive astronomical instruments, make it impossible for any observatory to completely escape the effects of RFI. There is an urgent need for strategies that will allow radio astronomers to share the spectrum with other users.

Ideally, one would subtract an interfering signal, leaving behind only the astronomical signal of interest and instrumental noise, with no loss of data. In practice, it is difficult to estimate and remove the interfering signal without biasing the underlying astronomical signal. It is, therefore, more common to identify and “flag” samples contaminated by RFI so that they can be ignored at some stage of processing, at the expense of losing a (potentially large) fraction of the data. The challenge then becomes robustly detecting RFI on short timescales, so as to maximize the fraction of usable data.

A number of RFI identification techniques have been developed. Some of these assume that signals from astrophysical sources can be

closely approximated as Gaussian random processes, calculate moments of the observed data, and flag non-Gaussian outliers as RFI (e.g. Nita et al. 2007; Nita & Gary 2010a; Purver et al. 2022). Others use principal component analysis to identify bases in which RFI stand out from sources of interest (Yuan et al. 2022). Machine learning offers another approach, in which algorithms are trained to recognize the same characteristics that humans use to manually identify RFI (e.g. Akeret et al. 2017; Vafaei Sadr et al. 2020; Pinchuk & Margot 2022). Each of these approaches has advantages and drawbacks. Statistical tests are straightforward and can be computationally inexpensive, but may also accidentally flag strong, impulsive astronomical sources. Principal component analysis and machine learning can use a rich, multi-dimensional representation of the data to identify RFI, but can fail when confronted with novel sources not in the training data set, though unsupervised learning methods may be able to overcome this weakness. Because RFI can take on many forms, and can have different impacts in different observing modes, it is important to explore new mitigation techniques that can complement and, in some cases, improve upon existing methods.

In this paper, we explore the use of cyclostationary signal processing (CSP) to identify RFI. A cyclostationary process is one with a statistical moment, such as mean or variance, that changes periodically or quasi-periodically (Gardner et al. 2006), as opposed to a wide-sense stationary process whose statistical moments are constant in time. Many sources of RFI are cyclostationary, with alternating current being a simple example. Cyclostationarity also arises from digital information encoding schemes in which the amplitude, frequency, and/or phase of a carrier wave switches between some finite number of possible states. Each state represents a symbol and the total num-

ber of possible states determines the number of bits that can be transmitted by each symbol. The signal will be cyclostationary at modulation frequencies related to the symbol rate, also known as the Baud rate. Since most¹ astrophysical processes are approximately wide-sense stationary, evidence of cyclostationarity could be a powerful way of distinguishing between RFI and astronomical sources.

Cyclostationarity has been discussed as an RFI mitigation technique in radio astronomy by Hellbourg et al. (2012) and Cucho-Padin et al. (2019), but has not yet been widely adopted. We have developed an algorithm for identifying and flagging RFI in astronomical data when there is significant evidence of cyclostationarity. Our long-term goal is to develop a system that can be integrated into modern radio astronomy digital spectrometers, but before doing so it is important that we determine the optimal algorithmic parameters and rigorously characterize its efficacy. As a first step in this process, we simulated a large number of human-generated signals using amplitude, phase, and frequency shift keying, and pre-processed them in a way that emulates the digital spectrometer used by the Robert C. Byrd Green Bank Telescope (GBT). Using this simulated data, we defined a “ground truth” that we then compared to the output of our algorithm. We simulated different symbol rates, numbers of bits per symbol, and signal-to-noise (S/N) ratios, in addition to the different keying techniques. This allowed us to explore the impact of different algorithmic parameters within a large parameter space. In §2 and §3 we provide some theoretical background and define our algorithm in detail. In §4 we describe our simulations, including the parameter space of the various signals and algorithmic parameters, and the metrics we use

¹ Pulsars and potential extraterrestrial techno-signatures are important exceptions.

to judge performance. We present results in §5 and discuss future avenues of research in §6, before concluding in §7.

2. OVERVIEW OF CYCLOSTATIONARY SIGNAL PROCESSING

Let $x(t)$ be a radio-frequency signal described by

$$x(t) = s(t)e^{2\pi if_c t + i\phi} + s^*(t)e^{-(2\pi if_c t + i\phi)} \quad (1)$$

where $s(t)$ is a signal of bandwidth B , $f_c \gg B$ is the carrier frequency, t is time, and ϕ is phase. $s(t)$ can itself be represented by in-phase and quadrature components:

$$s(t) = \frac{s_I(t) - is_Q(t)}{2} \quad (2)$$

If $s(t)$ is periodic on a timescale T_0 , then $x(t)$ will be cyclostationary, and we can extract several quantities of interest from $x(t)$. The first, known as the *non-conjugate cyclic autocorrelation function* (CAF), is a Fourier series representation of the traditional auto-correlation function given by

$$R_{xx^*}^\alpha(\tau) = \frac{1}{T_0} \int_{-T_0/2}^{T_0/2} E \left\{ x \left(t + \frac{\tau}{2} \right) x^* \left(t - \frac{\tau}{2} \right) \right\} e^{-2\pi i \alpha t} dt \quad (3)$$

where E is the expectation operator, $*$ denotes complex conjugation, t is time, τ is a time offset known as the lag, and α is the *cycle frequency* (Gardner 1991). A second quantity of interest

is the *conjugate* CAF²

$$R_{xx}^\alpha(\tau) = \frac{1}{T_0} \int_{-T_0/2}^{T_0/2} E \left\{ x \left(t + \frac{\tau}{2} \right) x \left(t - \frac{\tau}{2} \right) \right\} e^{-2\pi i \alpha t} dt \quad (4)$$

In a cyclostationary analog to the Wiener–Khinchin theorem, the Fourier transform of R_{xx}^α and R_{xx}^α with respect to τ yields the *non-conjugate* and *conjugate spectral correlation functions*, respectively (SCF; also known as the cyclic spectrum; Gardner 1991):

$$S_{xx^*}^\alpha(\nu) = \int_{-\infty}^{\infty} R_{xx^*}^\alpha(\tau) e^{-2\pi i \nu \tau} d\tau \quad (5)$$

$$S_{xx}^\alpha(\nu) = \int_{-\infty}^{\infty} R_{xx}^\alpha(\tau) e^{-2\pi i \nu \tau} d\tau. \quad (6)$$

We will refer to ν as the *spectral frequency* to differentiate it from the cycle frequency. The non-conjugate CAF and SCF will be non-zero only for $\alpha_n = n/T_0$, while the conjugate CAF and SCF will be non-zero only for $\alpha_n = n/T_0 \pm 2f_c$, where $n = 0, 1, 2, \dots$ is an integer. Note that when $\alpha = 0$, the non-conjugate SCF reduces to the usual definition of the power spectral density (PSD).

The non-conjugate and conjugate *spectral coherence* are normalized versions of the non-conjugate and conjugate SCF, defined as

$$\rho_{xx^*}^\alpha(\nu) = \frac{S_{xx^*}^\alpha(\nu)}{\sqrt{S_{xx^*}^0(\nu + \alpha/2) S_{xx^*}^0(\nu - \alpha/2)}} \quad (7)$$

$$\rho_{xx}^\alpha(\nu) = \frac{S_{xx}^\alpha(\nu)}{\sqrt{S_{xx^*}^0(\nu + \alpha/2) S_{xx^*}^0(\nu - \alpha/2)}} \quad (8)$$

² The nomenclature here can be confusing, since the non-conjugate CAF is calculated using the traditional definition of the autocorrelation function in which $x(t)$ is multiplied by a lagged version of its complex conjugate, while the conjugate CAF is calculated without using the conjugate of $x(t)$. We use this nomenclature to be consistent with other CSP literature.

where $S_{xx^*}^0(\nu \pm \alpha/2)$ is a frequency-shifted version of the PSD. Note that when $\alpha = 0$ the non-conjugate spectral coherence function is unity for all values of ν , regardless of the properties of the input signal.

Our algorithm exploits the fact that the SCF³ of a stationary process only has significant power when $\alpha = 0$, whereas the SCF of a cyclostationary process also has significant power at higher cycle frequencies. Since the magnitude of the spectral coherence function is ≤ 1 , it is especially useful for setting detection thresholds for data with arbitrary mean and variance.

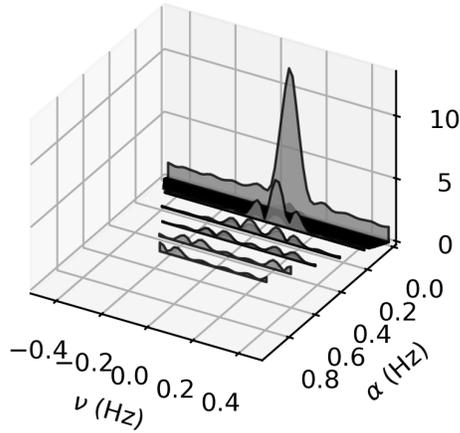
3. AN ALGORITHM FOR DETECTING RFI USING CYCLOSTATIONARY SIGNAL PROCESSING

In general, the data collected by a radio telescope may contain a large number of cyclostationary sources of RFI whose properties (e.g. carrier frequency, modulation frequency, encoding scheme, etc.) will not be known *a priori*. To blindly find evidence of cyclostationarity we need to have some way of efficiently estimating the SCF for a large number of discrete α . We make use of the *strip spectral correlation analyzer* (SSCA; Roberts et al. 1991), which works by time-averaging frequency-domain correlations (see Equations 9 and 10). Given a signal discretely sampled at a rate f_s with N total points, the SSCA estimates the SCF at N discrete values of α . The number of spectral frequencies, M , is controlled via a first-stage channelizer. In words, the steps in the SSCA are

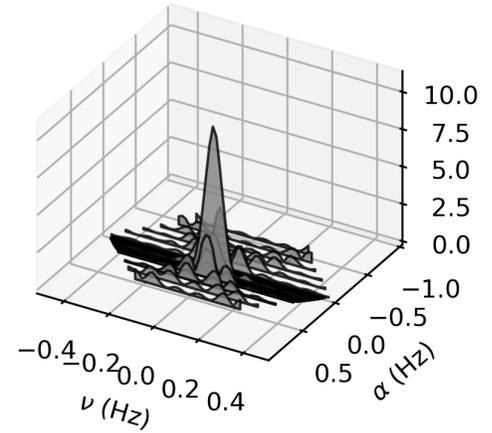
1. Take a data set, denoted as $x[n]$, of length N points and duration T .

³ In the remainder of this paper we will use SCF as an abbreviation for the non-conjugate and conjugate spectral correlation and coherence functions in contexts where these are interchangeable.

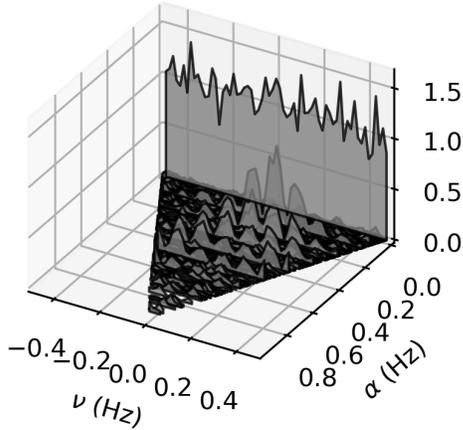
Non-Conjugate Spectral Correlation Function



Conjugate Spectral Correlation Function



Non-Conjugate Spectral Coherence Function



Conjugate Spectral Coherence Function

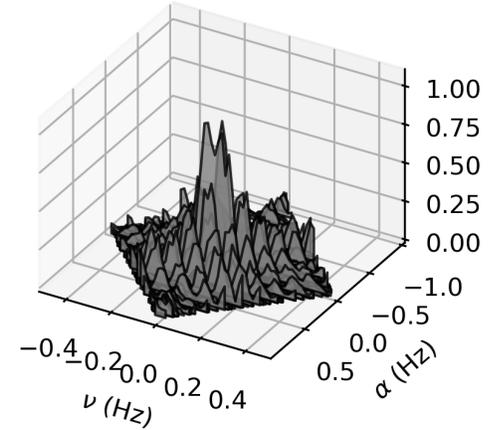


Figure 1. An example visualization of various forms of the SCF for a rectangular-pulse binary phase-shift keyed signal with a Baud rate of 0.1 Hz and carrier frequency of 0.05 Hz. The signal was 32,768 samples long and the SCFs were generated via our implementation of the strip spectral correlation analyzer using $M = 64$ (see text for details). For clarity, we have only plotted α corresponding to the top 200 values of the SCFs.

- | | |
|--|---|
| <p>230 2. Use a windowing function and sliding</p> <p>231 Fourier transform to channelize subsets of</p> <p>232 $x[n]$, each of length M, yielding $X[\nu_k, r]$,</p> <p>233 where ν_k are the channelizer frequencies</p> <p>234 (<i>not</i> the final spectral frequencies that ap-</p> <p>235 pear in Equations 5 and 6) and r is the</p> <p>236 time index.</p> | <p>240 4. Take a discrete Fourier transform of the</p> <p>241 result of step 3 along the time axis.</p> |
| <p>237 3. Multiply $X[\nu_k, r]$ by $x^*[r]$ (for the non-</p> <p>238 conjugate SCF) or $x[r]$ (for the conjugate</p> <p>239 SCF).</p> | <p>242 5. If desired, compute the spectral coher-</p> <p>243 ence using an over-sampled estimate of</p> <p>244 the PSD.</p> |

Table 1. Generalized Extreme Value Distribution Shape Parameters (ξ, μ, σ) for Various SSCA Parameters

Method	M						
	N	32	64	128	256	512	1024
Non-Conjugate	1024	(0.06,0.35,0.02)
	2048	(0.06,0.27,0.01)	(0.10,0.36,0.02)
	4096	(0.06,0.20,0.01)	(0.03,0.28,0.01)	(0.03,0.38,0.02)
	8192	(0.03,0.15,0.01)	(0.06,0.21,0.01)	(0.05,0.29,0.01)	(0.05,0.40,0.01)
	16384	(0.02,0.11,0.01)	(0.01,0.16,0.01)	(0.02,0.22,0.01)	(0.03,0.30,0.01)	(0.05,0.41,0.01)	...
	32768	(0.04,0.08,0.00)	(0.03,0.11,0.00)	(0.01,0.16,0.01)	(0.00,0.23,0.01)	(0.00,0.32,0.01)	(0.05,0.43,0.01)
Conjugate	1024	(0.06,0.42,0.03)
	2048	(0.00,0.32,0.02)	(0.06,0.43,0.02)
	4096	(0.02,0.25,0.01)	(0.06,0.33,0.02)	(0.03,0.44,0.02)
	8192	(0.00,0.18,0.01)	(0.00,0.25,0.01)	(0.04,0.34,0.02)	(0.06,0.45,0.02)
	16384	(0.06,0.14,0.01)	(0.00,0.19,0.01)	(0.00,0.26,0.01)	(0.00,0.35,0.02)	(0.04,0.46,0.02)	...
	32768	(0.00,0.10,0.01)	(0.02,0.14,0.01)	(0.02,0.19,0.01)	(0.00,0.27,0.01)	(0.05,0.36,0.02)	(0.03,0.48,0.02)

Mathematically, the SSCA for the non-conjugate SCF can be written as

$$\hat{S}_{xx^*}^{\nu_k+q\Delta\alpha} \left[n, \frac{\nu_k}{2} - \frac{q\Delta\alpha}{2} \right]_T = \sum_r X[\nu_k, r] x^*[r] w[n-r] e^{-2\pi i q r / N} \quad (9)$$

and for the conjugate SCF

$$\hat{S}_{xx}^{\nu_k+q\Delta\alpha} \left[n, \frac{\nu_k}{2} - \frac{q\Delta\alpha}{2} \right]_T = \sum_r X[\nu_k, r] x[r] w[n-r] e^{-2\pi i q r / N} \quad (10)$$

where the T subscript indicates time averaging, $\Delta\alpha = T^{-1}$ is the cycle frequency resolution, q is an integer index running from $-N/2$ to $N/2$, and w is a windowing function. The cycle and spectral frequencies are

$$\begin{aligned} \nu &= \frac{\nu_k}{2} - q \frac{\Delta\alpha}{2} \\ \alpha &= \nu_k + q\Delta\alpha. \end{aligned} \quad (11)$$

For the SSCA to provide an accurate estimate of the SCF it must satisfy the condition $N/M \gg 1$. In this work, we only considered cases where $N/M \geq 8$.

We implemented the SSCA in Python, closely following the approach described by Carter (1992). We use a Hann window and short-time

Fourier transform implemented as part of the cuSignal package (Thompson & Nicely 2021) for the first-stage channelization, overlapping each window by $M - 4$ samples. We use the CuPy (Okuta et al. 2017) fast Fourier transform routines for the second-stage transform. One critical and difficult aspect of using the SSCA to estimate the spectral coherence is estimating the PSD at the appropriate frequencies, particularly when the PSD estimate takes on small values, in which case small errors in the denominator of Eqs. 7 and 8 can lead to numerical artifacts. We use a time-averaged estimate of the PSD calculated from the input to the SSCA. Specifically, we use the SciPy implementation of Welch’s method with 32 time domain segments (i.e. each segment has a length of $N/32$ points), 50% overlap between segments, and a Hann window. Empirically, this leads to a robust estimate of the spectral coherence (see Fig. 1).

To use the SSCA to find evidence of cyclostationarity, we must define a robust detection statistic. We experimented with using the mean, median, and maximum energy of both $\hat{S}_{xx^{(*)}}^\alpha$ and $\hat{\rho}_{xx^{(*)}}^\alpha$ (in the remainder of this paper we use $(*)$ in the subscripts of S and ρ to mean both the non-conjugate and conjugate SCF). Recall that for all signal types, including stationary ones, the non-conjugate spectral

289 correlation function reduces to the PSD and
 290 the non-conjugate spectral coherence function
 291 reduces to unity for all ν . Since we are only
 292 interested in cyclostationary signals, it is there-
 293 fore sufficient to consider only $\alpha \neq 0$. We find
 294 that it is not ideal to use $\hat{S}_{xx^{(*)}}^\alpha$ for detection
 295 because the observed values depend on the in-
 296 put mean and variance of the data, which may
 297 not always be known in advance. We tried ac-
 298 counting for this by normalizing our input data
 299 to have zero mean and unit variance, but this
 300 biased $\hat{S}_{xx^{(*)}}^\alpha$ in the presence of strong signals.
 301 We had much better results using the $\hat{\rho}_{xx^{(*)}}^\alpha$ and
 302 so adopted this approach for all the results pre-
 303 sented here. Furthermore, as explained in §5.1,
 304 our algorithm works best when based on the
 305 maximum amplitude of the spectral coherence,
 306 as opposed to the mean or median.

307 The maximum amplitude of $\hat{\rho}_{xx^{(*)}}^\alpha$ follows a
 308 generalized extreme value (GEV) distribution,
 309 whose probability density function is

$$310 \quad f(x; \mu, \sigma, \xi) = \frac{1}{\sigma} t(x)^{\xi+1} e^{-t(x)} \quad (12)$$

311 where

$$312 \quad t(x) = \begin{cases} e^{-\frac{x-\mu}{\sigma}} & \text{if } \xi = 0, \\ [1 + \xi \left(\frac{x-\mu}{\sigma}\right)]^{-\frac{1}{\xi}} & \text{if } \xi \neq 0 \end{cases} \quad (13)$$

313 The quantile function for the GEV distribution
 314 is

$$315 \quad Q(p) = \begin{cases} \mu - \sigma \ln[-\ln(p)] & \text{if } \xi = 0, \\ \mu + \sigma \frac{[-\ln(p)]^{-\xi} - 1}{\xi} & \text{if } \xi \neq 0 \end{cases} \quad (14)$$

316 We can thus set a detection threshold, p_{thresh}
 317 such that we consider the data set under analy-
 318 sis to show significant evidence of cyclostation-
 319 arity when

$$320 \quad \max \left\{ |\rho_{xx^{(*)}}^{\alpha \neq 0}(\nu)| \right\}_{\text{observed}} > Q(p_{\text{thresh}}). \quad (15)$$

321 In principle the shape parameters should be in-
 322 dependent of the implementation details of the

323 SSCA, but in practice we find a small but com-
 324 plicated dependence on the choice of M and N ,
 325 and especially on the windowing function. We
 326 determined the shape parameters empirically
 327 for a Hann window for various combinations of
 328 M and N , and for the non-conjugate and conju-
 329 gate spectral coherence function, by generating
 330 normally distributed complex random values,
 331 passing the data through our SSCA implemen-
 332 tation, and recording $\max \left\{ |\rho_{xx^{(*)}}^{\alpha \neq 0}(\nu)| \right\}_{\text{observed}}$.
 333 We repeated this procedure 10^3 times and fit a
 334 GEV distribution to the results using the `stats`
 335 module in SciPy, recording the best-fit values
 336 of μ and β . The results are shown in Table
 337 1. Recall that we only considered cases where
 338 $M/N \geq 8$. We use these shape parameters to
 339 determine $Q(p_{\text{thresh}})$ for any given combination
 340 of M , N , and conjugate/non-conjugate spectral
 341 coherence.

342 We also explored using the mean and median
 343 values of the SCF as a detection statistic, which
 344 follow normal distributions in the presence of
 345 noise. Since they do not perform as well as the
 346 maximum value of the SCF, we do not report
 347 the distribution parameters here.

488 4. SIMULATIONS

349 We wish to measure the efficacy of our algo-
 350 rithm for various types of RFI and to determine
 351 the optimal values of M , N , and p_{thresh} . We
 352 are especially interested in emulating the data
 353 stream of modern radio telescope instruments
 354 so that our findings can be readily applied in
 355 real-world contexts. To do so, we simulated a
 356 large number of data sets and applied our al-
 357 gorithm for different parameter combinations.
 358 The steps in our simulations were

- 359 1. Define the signal parameters: symbol du-
 360 ration (t_{sym} , the inverse of the Baud rate),
 361 bits per symbol (n_{bit}), and energy per
 362 symbol (E_{sym}).

Table 2. Simulation and Algorithmic Parameters^a

Modulation Type	t_{sym} (samples)	n_{bit}	E_{sym}/N_0 (dB)	M	N	p_{thresh}
ASK	30	1	3	32	1024	0.001
OOK	32	2	5	64	2048	0.01
QAM	100	4	10	128	4096	0.05
FSK	128	6	20		8192	0.3
PSK	300	8			16384	0.6
	512				32768	0.9
	1000					0.95
	1024					0.99
						0.999
						0.9999

^aThis table is meant to be read down each column, and not across each row. Note that OOK signals are by definition limited to 1-bit.

- 363 2. Generate a symbol sequence, s , in the
364 form of random integers in the interval
365 $[0, 2^{n_{\text{bit}}})$, and use this to modulate some
366 property of a complex exponential carrier
367 wave. We simulated signals with seven
368 different types of modulation: amplitude
369 shift keying (ASK), on-off keying (OOK;
370 a special case of ASK), quadrature am-
371 plitude modulation (QAM; also a special
372 case of ASK), phase shift keying (PSK),
373 and frequency shift keying (FSK). .
 - 374 3. Add a simulated astrophysical spectral
375 line with a Gaussian profile.
 - 376 4. Include additive white Gaussian noise
377 (AWGN) with some noise power spectral
378 density (N_0).
 - 379 5. Pass the final time series through a simu-
380 lated astronomical spectrometer, produc-
381 ing a number of narrow-band, Nyquist-
382 sampled complex voltage time series cor-
383 responding to different frequency chan-
384 nels.
 - 385 6. Define a “ground truth” of which spec-
386 trometer channels and time samples con-
387 tain the simulated RFI.
 - 388 7. Independently analyze the output of each
389 spectrometer channel using our SSCA-
390 based algorithm using both the non-
391 conjugate and conjugate spectral coher-
392 ence function for various combinations of
393 M , N , and p_{thresh} .
 - 394 8. Compare the output of our algorithm with
395 the ground truth record and characterize
396 the performance of the algorithm using
397 various metrics.
 - 398 9. Repeat this process ten times for each sig-
399 nal parameter and algorithmic combina-
400 tion in order to better characterize the dis-
401 tribution of the various performance met-
402 rics.
- 403 In all cases we worked in normalized units, i.e.
404 with a sampling rate $f_s = 1$ Hz. The car-
405 rier frequency of the simulated RFI was $f_c =$
406 0.3 Hz. We always used a noise power of

407 $N_0 = 1 \text{ W Hz}^{-1}$, so that the S/N is equivalent
 408 to the value E_{sym} . The full parameter space of
 409 our simulations is shown in Table 2. In the fol-
 410 lowing sections we describe the above steps in
 411 more detail.

4.1. Simulated Interference Signals

413 There are many different modulation pro-
 414 cesses in use with telecommunications signals,
 415 some of which are quite complex. While we are
 416 interested in eventually characterizing our algo-
 417 rithm with as many encoding schemes as pos-
 418 sible, as a first step we limit our simulations
 419 to a simplified and somewhat idealized parame-
 420 ter space using basic amplitude, phase, and fre-
 421 quency shift keying processes. Each symbol se-
 422 quence, denoted as s , was a pulse train that
 423 consisted of n_{sym} symbols that were each t_{sym}
 424 length, so that s was a total of $n_{\text{sym}} \times t_{\text{sym}}$ sam-
 425 ples long. The symbols themselves were simply
 426 random integers in the interval $[0, 2^{n_{\text{bits}}}]$. This
 427 symbol sequence was convolved with a Hann
 428 window to reduce spectral leakage. The carrier
 429 wave for each signal was

$$430 \quad x(t) = \sqrt{\frac{E_{\text{sym}}}{t_{\text{sym}}}} e^{-2\pi i f_c t}. \quad (16)$$

431 Using this definition the integrated energy of
 432 $x(t)$ is $n_{\text{sym}} E_{\text{sym}}$. The modulation schemes are
 433 described below.

4.1.1. Amplitude Shift Keyed Signals

435 For generic ASK signals the modulated am-
 436 plitude is related to the symbol sequence by

$$437 \quad a = \frac{2s}{2^{n_{\text{bit}}} - 1} - 1. \quad (17)$$

438 This normalization ensures that the amplitude
 439 modulation is defined on the interval $[-1, +1]$.
 440 To ensure that the integrated energy is $n_{\text{sym}} E_{\text{sym}}$
 441 we divided the final signal by the standard de-
 442 viation of a . For the special case of an OOK

443 signal, $n_{\text{bit}} = 1$ and $a = s$ without any normal-
 444 ization, i.e. a is either 0 or 1.

445 We also simulated signals using QAM, which
 446 consists of two carrier waves, known as the
 447 in-phase (I) and quadrature (Q) components,
 448 which have the same frequency while being 90°
 449 out of phase. The amplitude of I and Q are
 450 modulated independently according to Equa-
 451 tion 17 using different symbol sequences. The
 452 total number of bits is split evenly between the
 453 two sequences. When $n_{\text{bit}} = 1$, $Q = 0$ and only
 454 I is used.

4.1.2. Phase Shift Keyed Signals

456 For PSK signals with $n_{\text{bit}} \geq 2$, the phase mod-
 457 ulation is given by

$$458 \quad \phi = \frac{2s + 1}{2^{n_{\text{bit}}}} \pi. \quad (18)$$

459 Using this definition the discrete phases are
 460 bounded on $[\pi/2^{n_{\text{bit}}}, \pi(2 - 1/2^{n_{\text{bit}}})]$. However,
 461 when $n_{\text{bit}} = 2$, we instead follow the typical
 462 convention that ϕ switches between 0 and π .

4.1.3. Frequency Shift Keyed Signals

464 We simulated a voltage controlled oscillator to
 465 generate FSK signals. The oscillator frequency
 466 was defined as

$$467 \quad f = f_0 + sK_0 \quad (19)$$

468 where f_0 is the quiescent oscillator frequency (in
 469 our case, the frequency of the carrier wave) and
 470 K_0 is the oscillator gain in units of Hz V^{-1} . We
 471 defined the phase of the carrier by integrating
 472 over f , thus ensuring that the phase was contin-
 473 uous across frequency shifts. In our simulations
 474 we used $K_0 = 0.01 \text{ Hz}$.

4.2. Simulated Spectral Line

476 Our algorithm should be insensitive to sta-
 477 tionary astronomical sources. We confirmed
 478 this by adding a voltage time series correspond-
 479 ing to a spectral line with a Gaussian profile.

480 In all our simulations the line had an ampli-
 481 tude of 20 V, was centered at a frequency of
 482 0.1 Hz, and had a full-width at half-maximum
 483 of 0.01 Hz. We first created the line with the rel-
 484 evant parameters in the frequency domain but
 485 with random phases, mimicking an incoherent
 486 astrophysical source. We then took an inverse
 487 Fourier transform to create the corresponding
 488 voltage time series.

489 4.3. Simulated Spectrometer

490 We passed the input data stream through a
 491 64-channel, 24-tap polyphase filterbank (PFB)
 492 spectrometer (Price 2021). This architecture
 493 is similar to that of the Versatile Green Bank
 494 Astronomical Spectrometer (VEGAS), the pri-
 495 mary backend for the GBT (Prestage et al.
 496 2015). In our implementation, we read $64 \times 24 =$
 497 1536 complex samples, multiplied this time se-
 498 ries by a windowing function of the same length,
 499 reshaped the data set into a 64×24 array, took
 500 a fast Fourier transform along the first axis,
 501 and then summed the result. This created an
 502 amplitude spectrum with 64 Nyquist-sampled
 503 channels. The window that we used was the
 504 product of a sinc function and Hann window.
 505 We did not form a power spectrum by taking
 506 the square modulus of the PFB output, but in-
 507 stead retained the full phase information. We
 508 repeated this channelization step until we accu-
 509 mulated $10N$ amplitude spectra.

510 4.4. Ground Truth Determination

511 The Hann window that we used to taper the
 512 symbol sequence and our PFB implementation
 513 both greatly reduce spectral leakage, but do not
 514 eliminate it completely. Therefore, the RFI sig-
 515 nal is present at some level across all PFB chan-
 516 nels, but usually at a level that is not expected
 517 to corrupt astronomical data. For the purposes
 518 of defining the ground truth comparison record,
 519 we passed both a noise-free version of the sig-
 520 nal and the realization of AWGN through our

521 PFB and formed the resulting power spectra.
 522 We considered the signal to be present at a sig-
 523 nificant level when its power was greater than
 524 or equal to the corresponding noise power.

525 4.5. SCF Estimation and Flagging

526 The output of the PFB was 64 narrow-band
 527 times series, each $10N$ points long. We ana-
 528 lyzed each channel independently in segments
 529 that were each N points long (recall that, in
 530 the SSCA, N is equal to the number of discrete
 531 α at which the SCF is estimated), resulting in
 532 ten SCF estimates for each PFB channel across
 533 our full data set. Note that there is a trade-off
 534 in the choice of N between cycle frequency res-
 535 olution and the time resolution with which we
 536 can flag data as being contaminated with RFI.

537 4.6. Performance Metrics

We computed several binary classification
 metrics. First, we compared the output of our
 algorithm for both the non-conjugate and con-
 jugate SCF to our ground truth definition and
 counted the number of true positives (TP), true
 negatives (TN), false positives (FP), and false
 negatives (FN). We also computed these for the
 union of the non-conjugate and conjugate out-
 puts. From these we calculated the following
 metrics:

$$\text{TPR} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (20)$$

$$\text{FPR} = \frac{\text{FP}}{\text{FP} + \text{TN}} \quad (21)$$

$$\phi = \frac{\text{TP} \times \text{TN} - \text{FP} \times \text{FN}}{\sqrt{(\text{TP} + \text{FP})(\text{TP} + \text{FN})(\text{TN} + \text{FP})(\text{TN} + \text{FN})}} \quad (22)$$

538 where TPR is the true positive rate, FPR is the
 539 false positive rate, and ϕ , also known as the
 540 Matthews correlation coefficient, is a widely-
 541 used binary classification metric that performs
 542 well for imbalanced classes. We also plot re-
 543 ceiver operating characteristic (ROC) curves,

ASK Combined Max

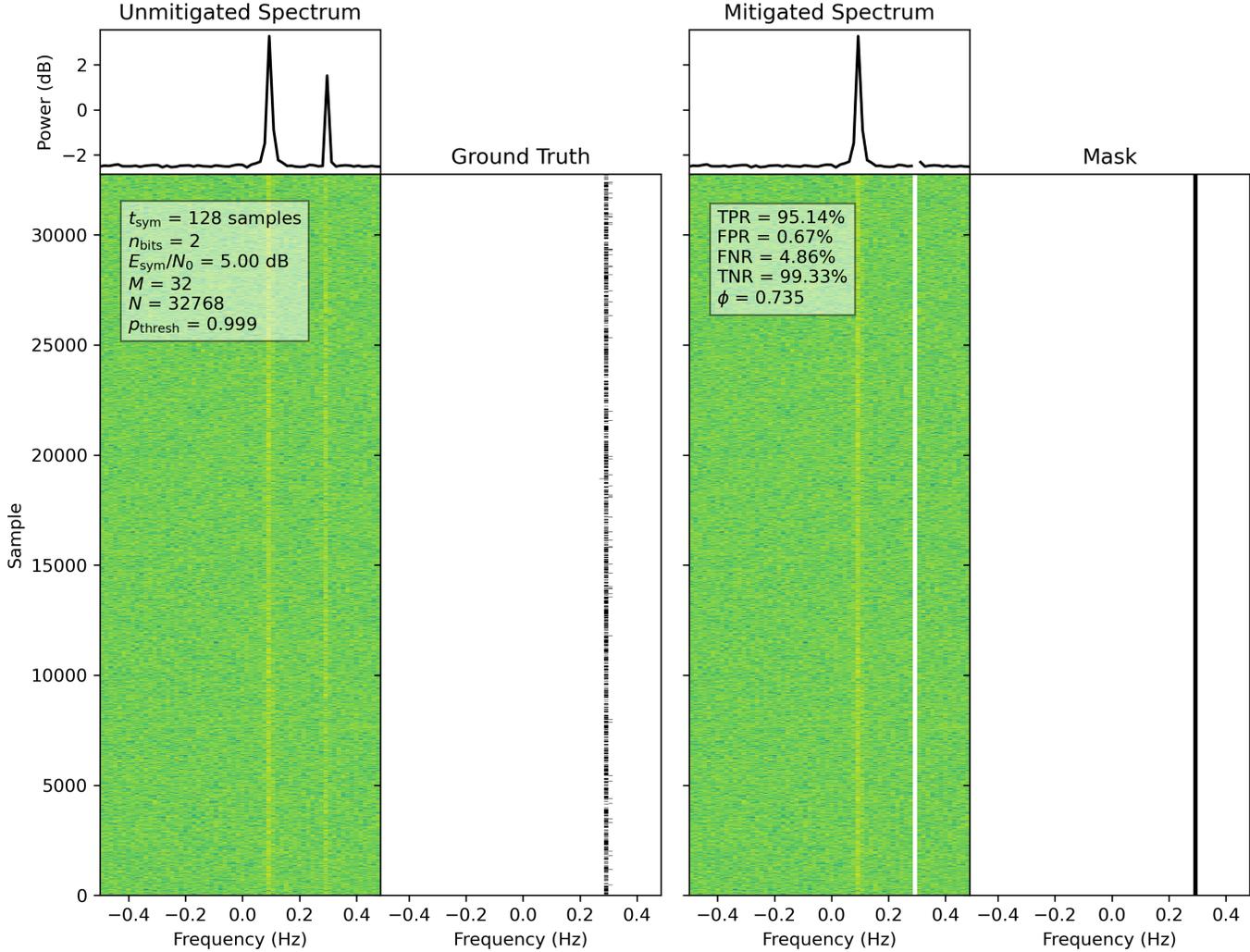


Figure 2. A summary plot for an ASK signal with $t_{\text{sym}} = 128$ samples, $n_{\text{bits}} = 2$, and $E_{\text{sym}}/N_0 = 5$ dB, processed using $(M, N) = (32, 32768)$ and $p_{\text{thresh}} = 0.999$. We show here the results of combining our algorithmic output for both the non-conjugate and conjugate SCF. This is one of the best performing combination of parameters in our simulations. In the mitigated spectrum we remove samples that are flagged by our algorithm (indicated in the mask panel), which completely removes the simulated signal. The simulated astrophysical spectral line appearing at 0.1 Hz is unaffected.

544 i.e. FPR vs TPR for different values of p_{thresh} ,
 545 and from these estimate the area-under-curve
 546 (AUC) value using a trapezoidal integration
 547 method as implemented in SciPy. A perfect
 548 classifier will have an ROC curve that imme-
 549 diately rises to a TPR of 1.0 and an FPR of
 550 0.0, and will maintain a TPR of 1.0 while the
 551 FPR rises as lower thresholds are used. The
 552 corresponding AUC would be 1.0. An uninfor-

553 mative classifier has an ROC curve with a slope
 554 of one and an AUC of 0.5. Values of ϕ and AUC
 555 in excess of 0.7 are generally considered to be
 556 good, and values in excess of 0.8 are generally
 557 considered to be very good.

5. RESULTS

559 Figure 2 shows an example summary plot from
 560 one of our simulations. We used 423,360 com-

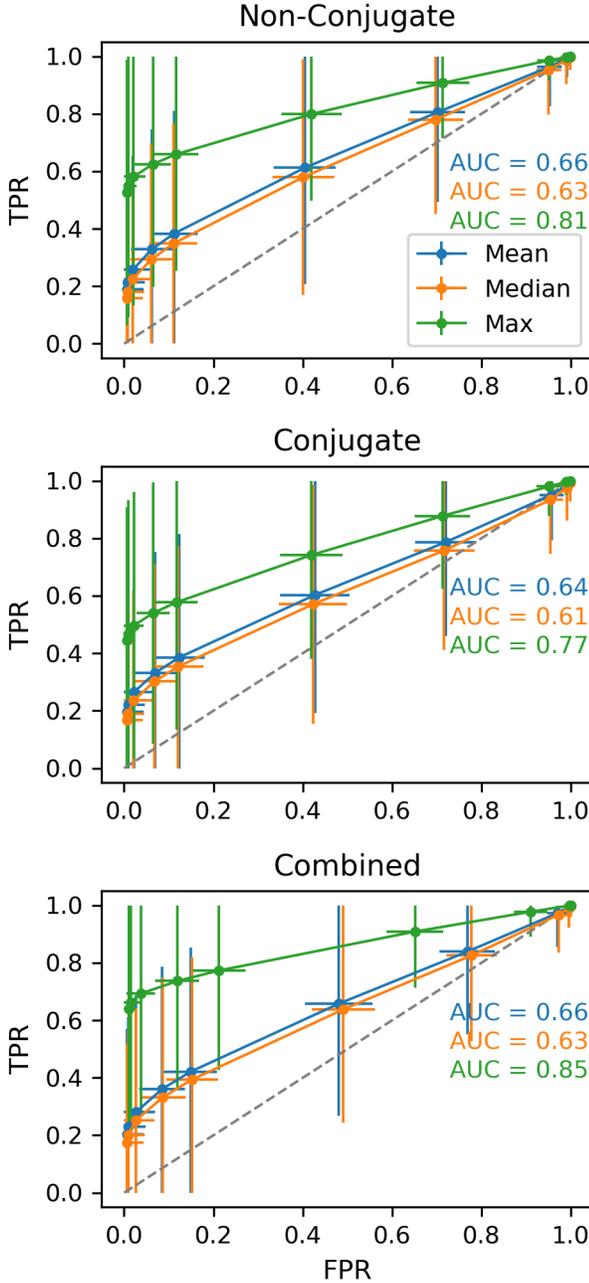


Figure 3. ROC curves with associated AUC values for different detection metrics using the non-conjugate (top) and conjugate (middle) SCF, and the combination of the two (bottom). The combined results using the maximum value of the SCF yields the highest AUC.

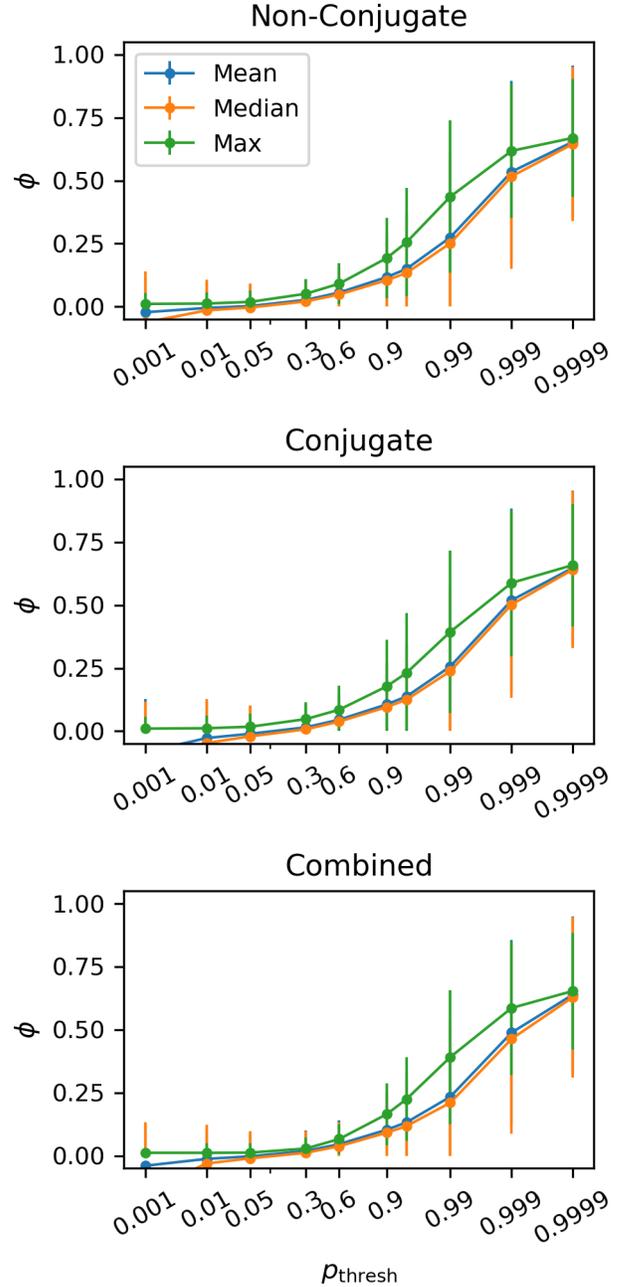


Figure 4. ϕ as a function of p_{thresh} for different detection metrics. Line colors are the same as in Fig. 3. The combined results using the maximum value of the SCF yields the highest ϕ at a value of 0.65 for $p_{\text{thresh}} = 0.9999$.

561 binations, and exploring this large parameter
 562 space is challenging. We begin by determin-
 563 ing whether it is most effective to flag based

564 on the mean, median, or maximum value of the
 565 SCF. Next, we find the optimal values of M , N ,
 566 and p_{thresh} , and then investigate how the perfor-
 567 mance varies with different signal properties.

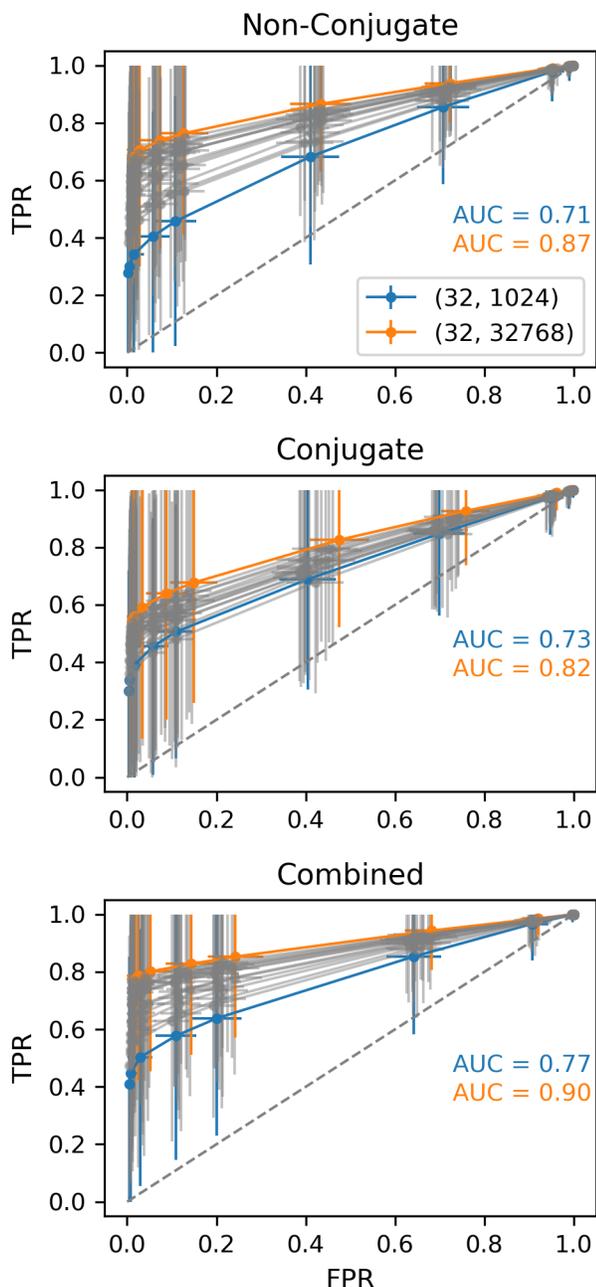


Figure 5. ROC curves for different combinations of M , and N . Curves for all combinations are plotted to provide a complete picture of the performance of our algorithm, but we highlight two combinations of interest. See text for details.

5.1. *Optimal Detection Metric*

Figures 3 and 4 show ROC curves and ϕ for different detection metrics when aggregating the

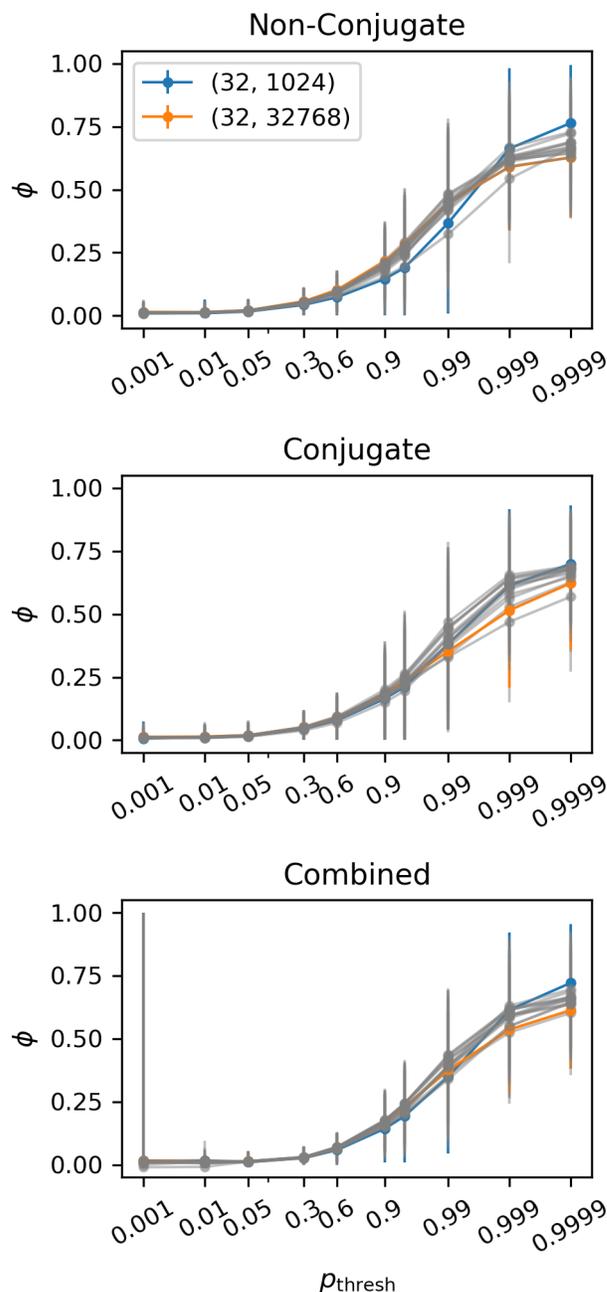


Figure 6. ϕ coefficients for different combinations of M , and N . As with Fig. 5, we plot all combinations but highlight two of interest.

results over all other simulation parameters. As noted previously, using the maximum value of the SCF significantly outperforms the mean or median, with an AUC of 0.85 and maximum ϕ value of 0.65 when combining results for the non-conjugate and conjugate SCF. We consider

577 these to be fairly good results, especially consid-
 578 ering that they cover a wide range of values of
 579 M and N , signal types, t_{sym} , n_{bits} , and E_{sym}/N_0 .
 580 In the remainder of this paper we will only con-
 581 sider results when using the maximum value of
 582 the SCF as a detection metric.

583 5.2. *Optimal (M, N) Pair*

584 In Figure 5 we show ROC curves and in Fig-
 585 ure 6 we show ϕ for all combinations of M and
 586 N . In both figures we averaged over all sig-
 587 nal properties (i.e., modulation type, t_{sym} , n_{bits} ,
 588 and E_{sym}/N_0). This provides the most complete
 589 measure of the performance of our algorithm
 590 but, as we will see, it includes signal properties
 591 for which the algorithm has weaknesses. We
 592 highlight two (M, N) pairs of particular inter-
 593 est, representing the highest AUC and ϕ .

594 The highest AUC is 0.90, which is obtained
 595 when using $(M, N) = (32, 32768)$ and the
 596 combination of the non-conjugate and conju-
 597 gate SCF. We consider this an excellent score.
 598 The lowest AUC is 0.71, which is obtained for
 599 $(M, N) = (32, 1024)$ when using only the non-
 600 conjugate SCF, which is still a good AUC score.
 601 However, the situation is reversed when consid-
 602 ering ϕ , i.e. the highest value is $\phi = 0.72$ for
 603 $(M, N) = (32, 1024)$ at $p_{\text{thresh}} = 0.9999$, while
 604 for $(M, N) = (32, 32768)$ and p_{thresh} , $\phi = 0.61$.

605 The discrepancy between AUC and ϕ can be
 606 understood by examining Table 3, which shows
 607 the TPR, FPR, TNR, and FNR for the two
 608 cases discussed above. As N increases from
 609 1024 to 32768, the TPR increases by a factor
 610 of 1.9, but the FPR increases by an even larger
 611 factor of 3.9. The ϕ coefficient punishes the al-
 612 gorithm for this larger relative increase in FPR.
 613 However, we note that the absolute improve-
 614 ment in TPR is 0.368, while the absolute de-
 615 terioration in FPR is only 0.0129, and remains
 616 quite low. In a real-world context, the question
 617 of which parameters are “better” will depend
 618 on the scientific goals of the observation. In

Table 3. Performance for $p_{\text{thresh}} = 0.9999$

(M, N)	TPR	FPR	TNR	FNR
(32, 1024)	0.409	0.00451	0.996	0.591
(32, 32768)	0.777	0.0174	0.983	0.223

619 some cases the large absolute improvement in
 620 TPR will make the slightly higher FPR tolera-
 621 ble, while other cases may require a lower FPR.
 622 For completeness we will report performance for
 623 both $(M, N) = (32, 1024)$ and $(32, 32768)$ in the
 624 remainder of this paper.

625 It is worth asking why the FPR increases with
 626 N ? From first principles, we would expect the
 627 SSCA to be more accurate as N increases be-
 628 cause it is a time-averaging technique for esti-
 629 mating the SCF, and by analyzing more data
 630 the signal-to-noise ratio of an interfering signal
 631 should go up. The observed behavior likely re-
 632 sults from our method for defining the ground
 633 truth comparison. Recall that we mark a sam-
 634 ple as truly containing RFI when the amplitude
 635 of the simulated signal is equal to the noise level.
 636 In the low signal-to-noise regime this will be sen-
 637 sitive to the exact realization of the noise. Such
 638 a situation can occur when RFI spills over with
 639 reduced amplitude into nearby PFB channels.
 640 Since we identify RFI in segments of length N ,
 641 the algorithm may flag data that technically
 642 falls just below the threshold for being included
 643 in our ground truth mask, leading to those sam-
 644 ples being marked as false positives. An anal-
 645 ogous situation could arise in the presence of
 646 transient RFI because good data will be flagged
 647 along with bad. We discuss potential ways to
 648 mitigate this shortcoming in §6.

649 We further note that, for any given value of
 650 N , there is a preference for smaller values of M ,

Table 4. Performance for Various p_{thresh}

N		0.898 ^a	0.95	0.984 ^b	0.99	0.998 ^c	0.9999	$1 - 4.61 \times 10^{-5}$ ^d
1024	TPR	0.641	0.577	0.519	0.502	0.456	0.409	0.239
	FPR	0.205	0.109	0.0410	0.0291	0.0100	0.00451	0.00254
32768	TPR	0.854	0.828	0.807	0.801	0.788	0.777	0.457
	FPR	0.247	0.143	0.0671	0.0526	0.0276	0.0174	0.0100

^aMinimizes Eq. 23 for $N = 1024$

^bMinimizes Eq. 23 for $N = 32768$

^cYields FPR = 0.01 for $N = 1024$

^dYields FPR = 0.01 for $N = 32768$

651 but the dependence on M is fairly weak. For the
 652 sake of simplicity we will only report results for
 653 $M = 32$ going forward. This is fortuitous be-
 654 cause smaller M reduce the computational com-
 655 plexity of the SSCA.

5.3. Optimal p_{thresh}

657 The optimal parameters for any classification
 658 algorithm will depend on the tolerance for false
 659 positives and false negatives, with ϕ being one
 660 commonly used measure. Figure 6 shows that
 661 ϕ is maximized for $p_{\text{thresh}} = 0.9999$ when aggre-
 662 gating over all simulation parameters, but there
 663 are diminishing returns for $p_{\text{thresh}} > 0.999$.

664 Another approach is to select a p_{thresh} based
 665 upon the ROC curves. A perfect classifier will
 666 always have TPR = 1 and FPR = 0, so we
 667 could choose the p_{thresh} that comes closest to
 668 this point. This is equivalent to finding the min-
 669 imum of

$$670 \sqrt{\text{FPR}^2(p_{\text{thresh}}) + [1 - \text{TPR}(p_{\text{thresh}})]^2}. \quad (23)$$

671 We used the SciPy `PchipInterpolator` rou-
 672 tine, which implements a piecewise cubic Her-
 673 mite interpolating polynomial, to interpolate
 674 FPR and TPR as a function of p_{thresh} , and then
 675 used Brent's method to find the minimum of
 676 Equation 23. For $N = 1024$ this results in

677 $p_{\text{thresh}}^{\text{opt}} = 0.898$ and for $N = 32768$ the opti-
 678 mal result is $p_{\text{thresh}}^{\text{opt}} = 0.984$ (recall that we only
 679 consider the case of $M = 32$).

680 Yet another approach is to choose a sensi-
 681 ble false alarm probability for the maximum
 682 value of the SCF to exceed what is expected
 683 for AWGN, e.g. $p_{\text{thresh}} = 0.95$ or 0.99 . While
 684 this may be attractive because it is motivated
 685 by the statistics for the SCF, we stress that it
 686 is *not* equivalent to the FPR of our algorithm,
 687 because we flag data in segments of length N
 688 (see §5.2 and §6).

689 In Table 4 we show TPRs and FPRs for vari-
 690 ous choices of p_{thresh} for our two representative
 691 values of N . For the remainder of this paper
 692 we will present results for $p_{\text{thresh}} = 0.99$ unless
 693 otherwise noted. We have chosen this value for
 694 three reasons: 1) it results in FPR $\lesssim 0.02$ for
 695 $N = 1024$ and FPR $\lesssim 0.05$ for $N = 32768$; 2) it
 696 is very close to the optimal value for $N = 32768$
 697 when using Eq. 23; and 3) it is one of the val-
 698 ues we directly simulated, avoiding the need to
 699 interpolate other results.

5.4. Performance for Different Modulation Types

702 In Figures 7 and 8 we show ROC curves and
 703 ϕ separated by modulation type for both $N =$
 704 1024 and $N = 32768$ (recall that we only con-

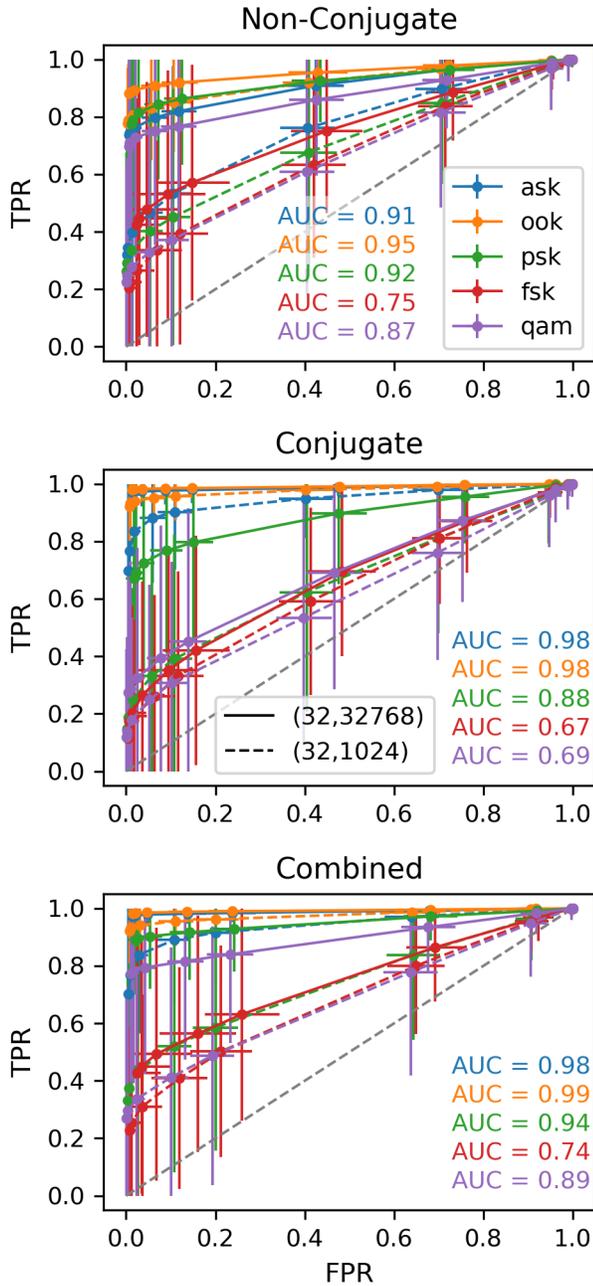


Figure 7. ROC curves and AUC values for different modulation types. Different color curves represent different modulation types, and different line styles indicate the two representative values of N that we consider. We see excellent performance for ASK, OOK, and PSK modulation, with good performance for QAM modulation, but a weakness to FSK signals.

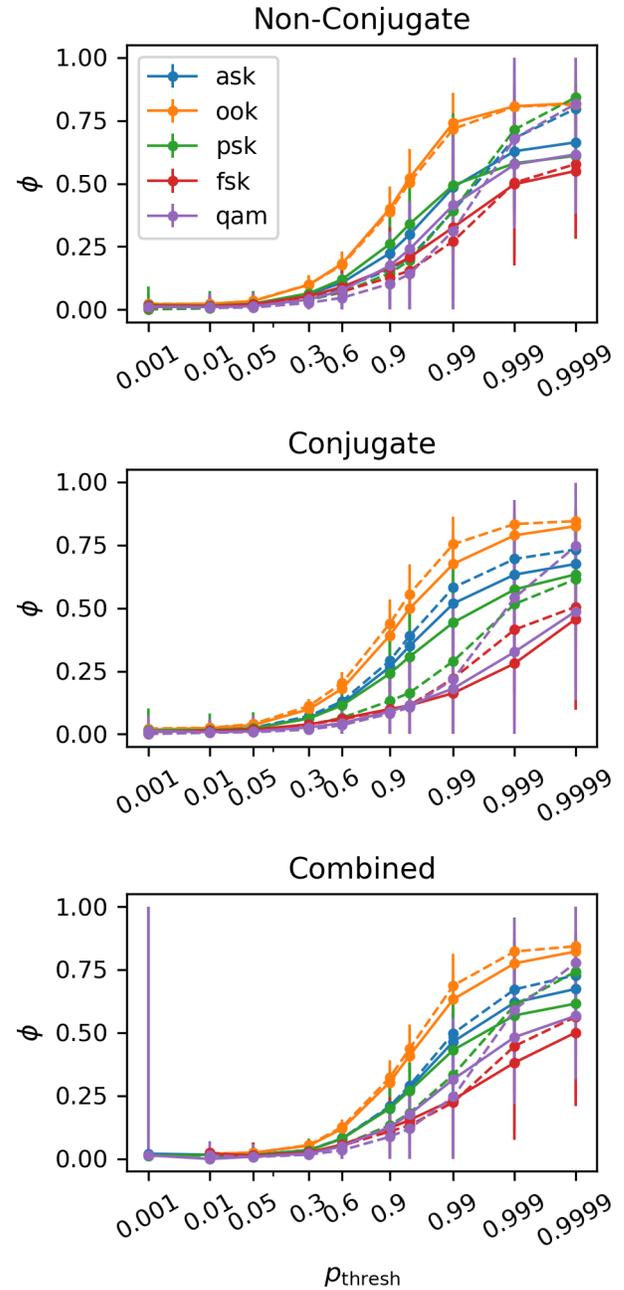


Figure 8. ϕ for different modulation types and values of N . The colors and line styles are the same as in Fig. 7. As with AUC, the results are best for ASK, OOK, and PSK modulation, and reasonably good for QAM modulation. However, the algorithm does not perform well for FSK modulation.

705 sider $M = 32$). Our algorithm works extremely

706 well for OOK signals, as well as ASK, PSK and
707 QAM signals, with maximum AUC $\gtrsim 0.9$. ϕ

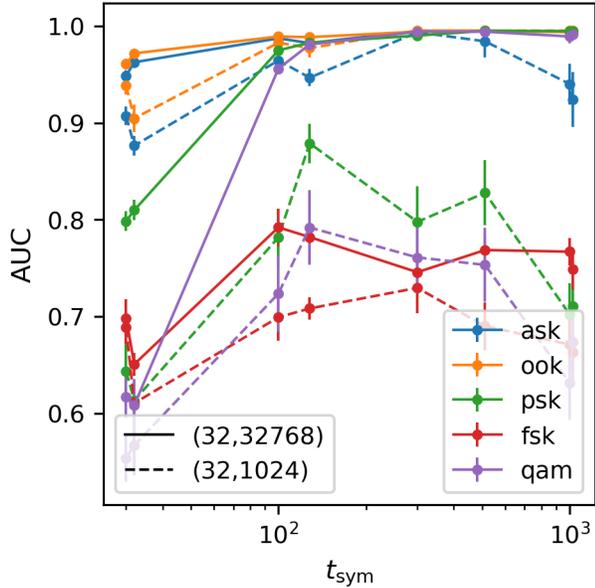


Figure 9. AUC values for different values of t_{sym} . Colors and line styles have the same meaning as in Fig. 7. AUC values increase quickly with larger t_{sym} , up to ~ 128 samples, and then either remain approximately constant or decrease slightly.

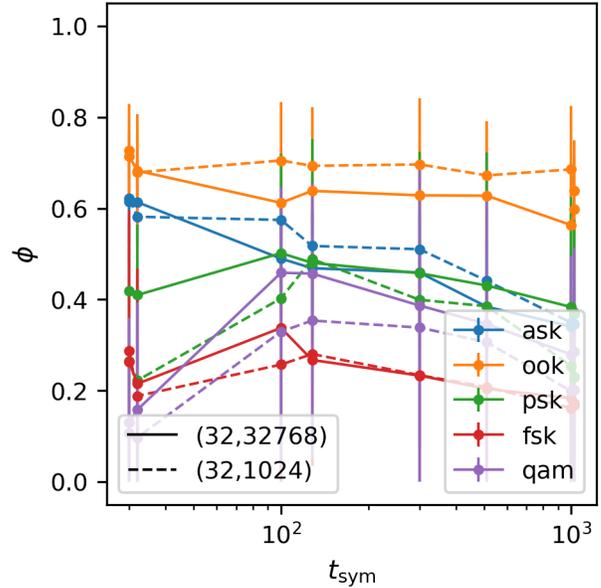


Figure 10. ϕ for different values of t_{sym} . Colors and line styles are the same as in Fig. 7. All values use $p_{\text{thresh}} = 0.99$. As with AUC values, ϕ increases quickly with larger t_{sym} , up to ~ 128 samples, and then either remain approximately constant or decrease slightly.

708 values span a wider range but are ≈ 0.75 for
 709 OOK signals at $p_{\text{thresh}} = 0.99$, and 0.6–0.75 for
 710 all other signal types (except FSK), albeit at
 711 higher p_{thresh} . The performance is not as good
 712 for FSK signals, with a maximum AUC $\simeq 0.74$
 713 and $\phi = 0.56$. As discussed in §6, the relative
 714 weakness to FSK signals most likely is a conse-
 715 quence of the way we independently analyze dif-
 716 ferent PFB channels. The frequency shift that
 717 is used could exceed the width of a PFB chan-
 718 nel, and in some cases the signal may not return
 719 to the original PFB channel within N samples,
 720 obscuring its cyclostationary nature. Neverthe-
 721 less, the performance for FSK signals is still rea-
 722 sonably good when aggregating over all other
 723 simulation parameters.

724 ASK and OOK signals are better detected us-
 725 ing the conjugate SCF, while PSK, FSK, and
 726 QAM signals are better detected using the non-
 727 conjugate SCF. The results for the combination
 728 of the two conjugation strategies are usually as

729 good, and in some cases slightly better, than the
 730 best individual results. This highlights the im-
 731 portance of using both conjugation strategies.
 732 For the sake of clarity, in the remainder of this
 733 paper we will only present the combined non-
 734 conjugate/conjugate results, but we will still
 735 separate results by modulation type since it has
 736 a significant impact on the performance of the
 737 the algorithm.

738 5.5. Performance for Different Symbol 739 Durations

740 In Figures 9 and 10 we show AUC values and
 741 ϕ for different values of t_{sym} , separated by mod-
 742 ulation type for our two representative combi-
 743 nations of N (we use $M = 32$ for both). ϕ
 744 coefficients are calculated for $p_{\text{thresh}} = 0.99$. As
 745 already noted in §5.4, the performance is best
 746 for OOK and ASK signals, followed by PSK
 747 and QAM, while performance is worst for FSK
 748 signals. However, we can now see that results

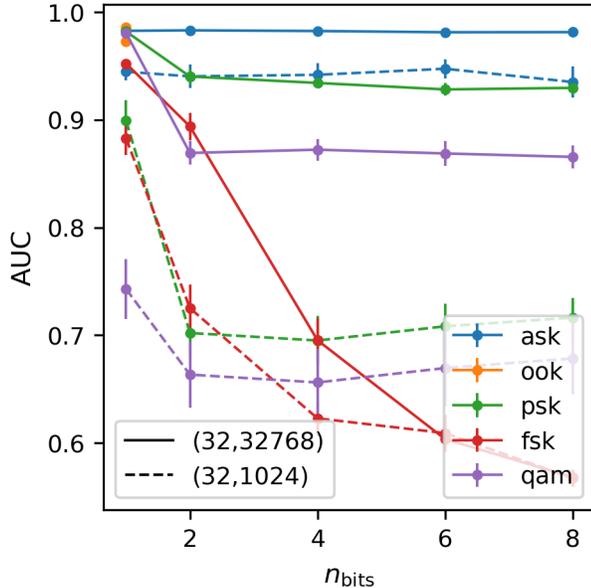


Figure 11. AUC values for different values of n_{bits} . Colors and line styles are the same as in Fig. 7. For most modulation types there is a small drop in performance between $n_{\text{bits}} = 1$ and 2, and relatively constant performance thereafter, though this does depend on the choice of N . However, the algorithm performs successively worse for FSK signals as n_{bits} increases, becoming nearly uninformative when $n_{\text{bits}} = 8$.

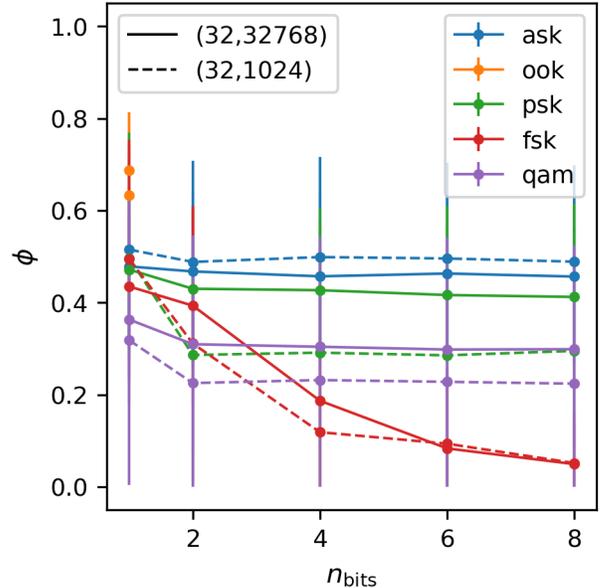


Figure 12. ϕ for different values of n_{bits} . Colors and line styles are the same as in Fig. 7. All values were calculated for $p_{\text{thresh}} = 0.99$. The drop in performance when going from $n_{\text{bits}} = 1$ to 2, is smaller than implied by AUC values. However, as with AUC values, ϕ coefficients imply that the algorithm is nearly uninformative when $n_{\text{bits}} = 8$.

749 also improve, sometimes significantly, when t_{sym}
750 increases from ~ 30 samples to ~ 100 sam-
751 ples, especially when measuring performance
752 via AUC. We can understand these trends by
753 recalling that we pass our data through a first-
754 stage 64-channel PFB, and analyze each chan-
755 nel independently. When $t_{\text{sym}} \lesssim 64$, signals are
756 spread across multiple PFB channels, reducing
757 the signal-to-noise ratio. Once t_{sym} is greater
758 than the width of a PFB channel, the signal is
759 fully contained within one channel and the per-
760 formance of the algorithm does not change very
761 much, until reaching the highest values of t_{sym} .

762 At the highest values of t_{sym} we do see a signifi-
763 cant drop in performance when using $N = 1024$,
764 because as t_{sym} increases there are fewer sym-
765 bols over which we can average to obtain an ac-
766 curate estimate of the SCF. This is an argument

767 against using small values of N when trying to
768 detect narrow-bandwidth signals.

769 The algorithm continues to perform poorly for
770 FSK signals because the frequency shift can still
771 exceed the width of a PFB channel.

772 We tested t_{sym} that are and are not evenly di-
773 visible by N , i.e. that have or do not have Baud
774 rates that align precisely with the SSCA cycle
775 frequency bins (see Table 2 for the complete list
776 of t_{sym}). As expected, the performance for Baud
777 rates that are not equal to a cycle frequency
778 bin are somewhat lower than similar Baud rates
779 that do align with a cycle frequency bin. We re-
780 turn to this point in §6.

781 5.6. Performance for Different Numbers of 782 Bits

783 Figures 11 and 12 show AUC values and ϕ as
784 function of n_{bits} per symbol, separated by mod-

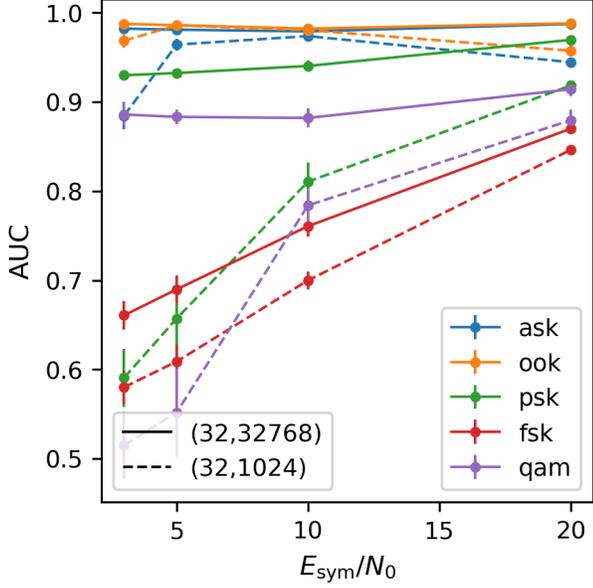


Figure 13. AUC values for different E_{sym}/N_0 . Colors and line styles are the same as in Fig. 9. As expected, the algorithm performs better as E_{sym}/N_0 increases for all modulation types.

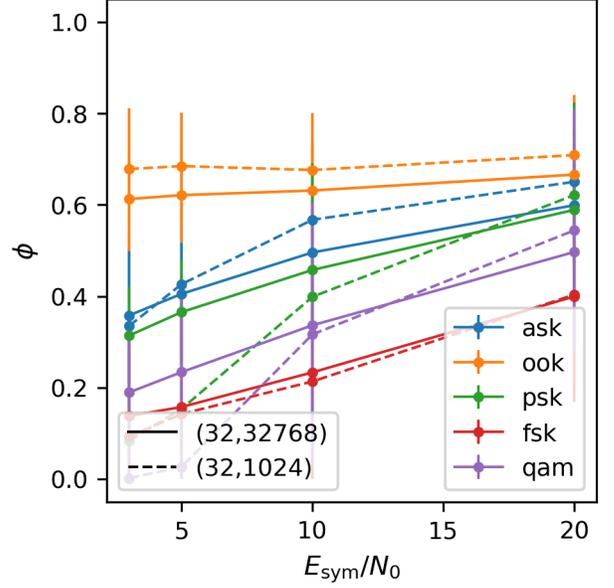


Figure 14. ϕ for different E_{sym}/N_0 . Colors and line styles are the same as in Fig. 9. All values were calculated for $p_{\text{thresh}} = 0.99$. As with AUC values, the ϕ coefficients show that the algorithm performs better as E_{sym}/N_0 increases for all modulation types.

785 modulation type, for our two representative combi-
 786 nations of N (both using $M = 32$). ϕ coeffi-
 787 cients are calculated for $p_{\text{thresh}} = 0.99$. There is
 788 a slight drop in performance when going from
 789 one bit to two for ASK, PSK, and QAM signals,
 790 but no significant dependence on n_{bits} at higher
 791 values (OOK signals are only 1-bit). However,
 792 there is a strong dependence on n_{bits} for FSK
 793 signals, with more bits per symbol leading to
 794 steadily worse performance. Once again, this
 795 is related to our approach of analyzing PFB
 796 channels independently. In our implementation,
 797 FSK-like signals with more bits per symbol will
 798 be spread over a wider range of frequencies. In
 799 the extreme case, a signal may not return to
 800 a given frequency channel within N samples, in
 801 which case its cyclostationary nature will not be
 802 detected at all by our algorithm. This does in-
 803 deed seem to be the case, as can be seen by AUC
 804 values approach 0.5 and ϕ coefficients approach
 805 zero as n_{bits} increases. However, we can also
 806 see that the algorithm performs well for FSK

807 signals when $n_{\text{bits}} = 1$, and its performance re-
 808 mains acceptable up to $n_{\text{bits}} = 2-4$.

5.7. Performance for Different E_{sym}/N_0

810 Figure 14 shows ϕ as a function of E_{sym}/N_0 for
 811 different modulation types, for both $N = 1024$
 812 and $N = 32768$, and using $p_{\text{thresh}} = 0.99$. As
 813 expected, higher E_{sym}/N_0 leads to better per-
 814 formance. The relative improvement is not as
 815 high for ASK, OOK, PSK, and QAM signals
 816 since the algorithm already detects these signals
 817 well, even at low E_{sym}/N_0 , but there is a large
 818 relative improvement for FSK signals. However,
 819 as discussed in §5.5 and 5.6, there is a strong de-
 820 pendence on other parameters for FSK signals.
 821 The improvements seen here are due to those
 822 few cases where the algorithm works reasonably
 823 well for FSK signals (e.g. $n_{\text{bits}} = 1$). For others,
 824 such as very high values of n_{bits} , the algorithm
 825 does not work well for FSK signals even at very
 826 high E_{sym}/N_0 .

5.8. Performance for Simulated Spectral Line

As noted above, our algorithm should not identify astrophysical spectral lines as potential sources of RFI because they are not cyclostationary. All of the results discussed in the preceding sections include a simulated spectral line in the data, and so any false positives would include samples containing signal from this line. To further verify that this simulated line is not mistakenly being identified as RFI, we also simulated data sets containing only the line and analyzed them using all algorithmic parameter combinations and recorded the FPR (since there is no true source of RFI, the TPR is undefined). We did the same for pure AWGN. We then performed a two-sided Kolmogorov-Smirnov test using SciPy’s `kstest` routine. We find KS test statistics of 0.0078, 0.0064, and 0.0096 for the non-conjugate SCF, conjugate SCF, and combined results, respectively. These correspond to p -values of ≥ 0.99 . As expected, we thus find no evidence for rejecting the null hypothesis that the FPRs of the data sets containing the simulated spectral line and pure noise come from the same distribution.

6. DISCUSSION

These results show that cyclostationary tests are a promising approach to RFI mitigation. Aggregating our results across different signal properties provides a more complete picture of how our algorithm performs, but for particularly favorable combinations of signal properties the performance can be much better than the aggregate results imply. As an example, when using the combined non-conjugate/conjugate SCF, $(M, N) = (32, 32768)$ and $p_{\text{thresh}} = 0.99$, for $t_{\text{sym}} = 128$ samples, $n_{\text{bits}} = 1$, and $E_{\text{sym}}/N_0 = 10$ dB, we find $\text{TPR} > 0.97$ and $\text{FPR} < 0.06$ for all modulation types except FSK (which has $\text{TPR} = 0.94$ and $\text{FPR} = 0.13$). If we choose $(M, N) = (32, 1024)$ and $p_{\text{thresh}} = 0.999$ we can achieve $\phi > 0.78$

for all modulation types except QAM ($\phi = 0.72$) and FSK ($\phi = 0.69$). Furthermore, we find no evidence that our algorithm systematically flags the simulated astrophysical spectral line that we included in our simulations. Obviously, we cannot optimize the properties of real-world RFI to maximize the effectiveness of mitigation techniques, but these results do suggest that cyclostationary tests can perform extremely well and potentially complement other approaches. For example, spectral Kurtosis (Nita & Gary 2010a,b; Smith et al. 2022) is a computationally simple statistical method that distinguishes normally-distributed data and from RFI, though it has weaknesses to sidelobe spillover as well as weaker signals and those that have a 50% duty cycle. Smith et al. (2022) measured the performance of single- and multi-scale SK using many of the same simulated sources of RFI as we use here. ϕ scores varied substantially depending on the characteristics of the signal and SK parameters, but could be as high as ~ 0.75 for ASK signals with high data rates, and were usually ~ 0.5 – 0.7 , which is broadly similar to our results (E. Smith, private communication). AOFLAGGER is used on low-frequency arrays such as the Low Frequency Array (LOFAR) and the Murchison Widefield Array (MWA; Offringa et al. 2010a,b, 2012), and flags the post-correlation visibilities with the highly optimized `SumThreshold` method. AOFLAGGER does very well at flagging most RFI in the dataset, but operates on the power values, which means its performance may be hindered by uneven bandpass responses or strong periodic astronomical signals such as pulsars or FRBs.

Nevertheless, the current implementation of our algorithm does have a weakness to signals that are spread across multiple PFB channels, whether because of the modulation technique being used or intrinsic bandwidth of the signal itself. As noted previously, this stems from ana-

912 lyzing each PFB channel independently. When
 913 a signal is spread across multiple channels the
 914 effective E_{sym}/N_0 decreases. In the worst-case
 915 scenario, a frequency-switched signal may not
 916 return to a given PFB channel within the block
 917 of data that we analyze, completely obscur-
 918 ing its cyclostationary nature. We chose to
 919 analyze PFB channels independently in order
 920 to more closely match the architecture of the
 921 digital spectrometer used at the GBT. In this
 922 case, PFB channels are formed on a field pro-
 923 grammable gate array prior to being transmit-
 924 ted to computers where a real-time RFI miti-
 925 gation algorithm might be implemented. How-
 926 ever, there are alternative architectures or ap-
 927 proaches. For example, RFI mitigation could
 928 be implemented prior to the PFB. The PFB op-
 929 eration could also be inverted in software, and
 930 groups of PFB channels could be analyzed in
 931 groups covering sufficient bandwidth to capture
 932 even relatively broad-band RFI. These groups
 933 could also be made to overlap by using an over-
 934 sampled PFB, to avoid missing signals that
 935 cross over group boundaries.

936 In §5.2 we showed that while using larger val-
 937 ues of N leads to a higher TPR, it also leads to
 938 a higher FPR. We attribute this to our method
 939 of defining a ground truth comparison, but we
 940 also expect it to be true when data contain tran-
 941 sient RFI. Our algorithm operates on data in
 942 segments of length N , so there is a chance that
 943 samples that are free of RFI will be incorrectly
 944 flagged when RFI only contaminates some of
 945 the samples. By analyzing data that contains
 946 both cyclostationary and non-cyclostationary
 947 signals, we would also lower the sensitivity of
 948 the algorithm. We could avoid these pitfalls by
 949 using multiple values of N and selecting the best
 950 value for any given segment of data by choos-
 951 ing the value that maximizes the signal-to-noise
 952 ratio of the SCF. Using multiple values of N
 953 would also lead to different cycle frequency res-
 954 olutions, which could help detect signals at dif-

955 ferent Baud rates. However, this would increase
 956 computational cost, which is already high to be-
 957 gin with (see below).

958 Quantization errors may also impact the per-
 959 formance of our algorithm. In our simulations
 960 we generated signals with floating point pre-
 961 cision, but modern analog-to-digital converters
 962 (ADCs) use much lower quantization depth, e.g.
 963 the VEGAS spectrometer used at the GBT out-
 964 puts 8-bit values. This may be at least partially
 965 ameliorated by using higher bit-depth ADCs —
 966 commercial models are now available that out-
 967 put 12-bit values and that can sample band-
 968 widths of several GHz.

969 However, each of these approaches does come
 970 with challenges. Sampling with more bits
 971 increases data rates, requiring new network
 972 topologies. Analyzing the full bandwidth with
 973 different values of N would be computationally
 974 expensive and may exceed the resources avail-
 975 able with modern hardware for all but rela-
 976 tively narrow observing bands. Inverting the
 977 PFB operation and using overlapping groups
 978 also adds computational cost. The compu-
 979 tational complexity of the SSCA algorithm is
 980 $O \sim NM \log_2 N$ (Roberts et al. 1991). Process-
 981 ing a bandwidth of 1 GHz in two separate po-
 982 larization channels with our optimal algorithmic
 983 parameters of $M = 32$ and $N = 32768$ in real-
 984 time would thus require a computing system ca-
 985 pable of ~ 30 PFLOPS. This is well beyond the
 986 capability of current commercial graphics card,
 987 but the current generation of GPUs designed for
 988 artificial intelligence training offer theoretical
 989 maximum computational power of $\sim 300\text{--}600$
 990 TFLOPS, depending on the numerical precision
 991 being used. Over the next several years it may
 992 become feasible to adopt a hybrid approach,
 993 wherein wide observing bandwidths are split
 994 into a modest number of overlapping sub-bands
 995 and processed independently before being com-
 996 bined to record the full bandwidth. A similar
 997 approach has already been developed to enable

998 real-time coherent dedispersion of pulsars using
 999 the GBT’s 0.7 – 4 GHz ultrawideband receiver,
 1000 for which the GBT’s VEGAS spectrometer will
 1001 use 24 compute nodes to process 3.3 GHz of
 1002 instantaneous bandwidth, as well as to enable
 1003 cyclostationary techniques for studying pulsars
 1004 (Demorest 2011). Another approach is to es-
 1005 chew real-time RFI mitigation in favor of tem-
 1006 porarily recording Nyquist-sampled voltages to
 1007 disk and processing them offline with some rea-
 1008 sonable turnaround time. This approach is used
 1009 by the Breakthrough Listen project (MacMahon
 1010 et al. 2018) to process several GHz of instanta-
 1011 neous bandwidth. We leave a detailed analysis
 1012 of these approaches to future work.

1013 We chose a limited number of idealized signal
 1014 types to illustrate a CSP-based approach to RFI
 1015 mitigation, but real-world telecommunications
 1016 signals can be much more complex. In future
 1017 work we plan on simulating additional modula-
 1018 tion strategies and windowing functions, includ-
 1019 ing more complex astrophysical sources, and
 1020 adding multiple sources of RFI within the fre-
 1021 quency range of interest. More complex strate-
 1022 gies for improving our algorithm could also in-
 1023 clude using multiple PFBs to channelize the
 1024 data with different numbers of channels. Fi-
 1025 nally, as an alternative to our blind identifica-
 1026 tion algorithm, we could study the local RFI
 1027 environment and use cyclostationary detectors
 1028 that are tuned to sources of RFI with known
 1029 properties, which would greatly reduce the com-
 1030 putational cost. We also plan to apply our al-
 1031 gorithm using the optimal parameters derived

1032 here to archived astronomical data collected
 1033 with the GBT.

1034 7. CONCLUSIONS

1035 We have developed an approach to identify-
 1036 ing and mitigating RFI by testing whether data
 1037 contain significant evidence of cyclostationar-
 1038 ity, and tested its performance using a range of
 1039 simulated signals. We find good performance
 1040 for most simulated signals, with some weak-
 1041 nesses to broad-band and frequency-switched
 1042 signals. Specifically, when using optimal algo-
 1043 rithmic parameters we find AUC scores > 0.90
 1044 and ϕ scores $\gtrsim 0.61$, aggregated over all mod-
 1045 ulation schemes, symbol durations, bits-per-
 1046 symbol, and signal-to-noise ratios that we sim-
 1047 ulated. The algorithm performs best for OOK
 1048 signals and reasonably well for more generic
 1049 ASK and PSK signals. We find no systemic ten-
 1050 dency for our algorithm to incorrectly identify
 1051 a simulated astrophysical spectral line. We be-
 1052 lieve that tests of cyclostationarity are a promis-
 1053 ing technique for RFI mitigation that can com-
 1054 plement other approaches.

1055 This work is supported by the National Sci-
 1056 ence Foundation through Advanced Technolo-
 1057 gies and Instrumentation grant #1910302. We
 1058 are grateful to an anonymous referee for provid-
 1059 ing comments that improved the quality of this
 1060 manuscript, and to Chad Spooner for helpful
 1061 discussions and for maintaining [cyclotatory.](#)
 1062 [blog](#).

REFERENCES

- 1063 Akeret, J., Chang, C., Lucchi, A., & Refregier, A.
 1064 2017, *Astronomy and Computing*, 18, 35,
 1065 doi: [10.1016/j.ascom.2017.01.002](https://doi.org/10.1016/j.ascom.2017.01.002)
 1066 Carter, N. J. 1992, Master’s thesis, Naval
 1067 Postgraduate School, Monterey, CA.
 1068 Cucho-Padin, G., Wang, Y., Li, E., et al. 2019,
 1069 *Radio Science*, 54, 986,
 1070 doi: [10.1029/2019RS006902](https://doi.org/10.1029/2019RS006902)
- 1071 Demorest, P. B. 2011, *MNRAS*, 416, 2821,
 1072 doi: [10.1111/j.1365-2966.2011.19230.x](https://doi.org/10.1111/j.1365-2966.2011.19230.x)
 1073 Gardner, W. A. 1991, *IEEE Signal Processing*
 1074 *Magazine*, 8, 14, doi: [10.1109/79.81007](https://doi.org/10.1109/79.81007)
 1075 Gardner, W. A., Napolitano, A., & Paura, L.
 1076 2006, *Signal Processing*, 86, 639, doi: <https://doi.org/10.1016/j.sigpro.2005.06.016>
 1077

- 1078 Hellbourg, G., Weber, R., Capdessus, C., &
1079 Boonstra, A.-J. 2012, *Comptes Rendus*
1080 *Physique*, 13, 71, doi: [10.1016/j.crhy.2011.10.010](https://doi.org/10.1016/j.crhy.2011.10.010)
1081 MacMahon, D. H. E., Price, D. C., Lebofsky, M.,
1082 et al. 2018, *PASP*, 130, 044502,
1083 doi: [10.1088/1538-3873/aa80d2](https://doi.org/10.1088/1538-3873/aa80d2)
1084 Nita, G. M., & Gary, D. E. 2010a, *MNRAS*, 406,
1085 L60, doi: [10.1111/j.1745-3933.2010.00882.x](https://doi.org/10.1111/j.1745-3933.2010.00882.x)
1086 —. 2010b, *PASP*, 122, 595, doi: [10.1086/652409](https://doi.org/10.1086/652409)
1087 Nita, G. M., Gary, D. E., Liu, Z., Hurford, G. J.,
1088 & White, S. M. 2007, *PASP*, 119, 805,
1089 doi: [10.1086/520938](https://doi.org/10.1086/520938)
1090 Offringa, A. R., de Bruyn, A. G., Biehl, M., et al.
1091 2010a, *MNRAS*, 405, 155,
1092 doi: [10.1111/j.1365-2966.2010.16471.x](https://doi.org/10.1111/j.1365-2966.2010.16471.x)
1093 Offringa, A. R., de Bruyn, A. G., Zaroubi, S., &
1094 Biehl, M. 2010b, arXiv e-prints,
1095 arXiv:1007.2089.
1096 <https://arxiv.org/abs/1007.2089>
1097 Offringa, A. R., van de Gronde, J. J., & Roerdink,
1098 J. B. T. M. 2012, *A&A*, 539, A95,
1099 doi: [10.1051/0004-6361/201118497](https://doi.org/10.1051/0004-6361/201118497)
1100 Okuta, R., Unno, Y., Nishino, D., Hido, S., &
1101 Loomis, C. 2017, in *Proceedings of Workshop*
1102 *on Machine Learning Systems (LearningSys) in*
1103 *The Thirty-first Annual Conference on Neural*
1104 *Information Processing Systems (NIPS)*.
1105 [http://learningsys.org/nips17/assets/papers/](http://learningsys.org/nips17/assets/papers/paper_16.pdf)
1106 [paper_16.pdf](http://learningsys.org/nips17/assets/papers/paper_16.pdf)
1107 Pinchuk, P., & Margot, J.-L. 2022, *AJ*, 163, 76,
1108 doi: [10.3847/1538-3881/ac426f](https://doi.org/10.3847/1538-3881/ac426f)
1109 Prestage, R. M., Bloss, M., Brandt, J., et al. 2015,
1110 in *2015 URSI-USNC Radio Science Meeting*, 4,
1111 doi: [10.1109/USNC-URSI.2015.7303578](https://doi.org/10.1109/USNC-URSI.2015.7303578)
1112 Price, D. C. 2021, in *The WSPC Handbook of*
1113 *Astronomical Instrumentation, Volume 1: Radio*
1114 *Astronomic al Instrumentation*, ed.
1115 A. Wolszczan, 159–179,
1116 doi: [10.1142/9789811203770_0007](https://doi.org/10.1142/9789811203770_0007)
1117 Purver, M., Bassa, C. G., Cognard, I., et al. 2022,
1118 *MNRAS*, 510, 1597,
1119 doi: [10.1093/mnras/stab3434](https://doi.org/10.1093/mnras/stab3434)
1120 Roberts, R. S., Brown, W. A., & Loomis,
1121 Herschel H., J. 1991, *IEEE Signal Processing*
1122 *Magazine*, 8, 38, doi: [10.1109/79.81008](https://doi.org/10.1109/79.81008)
1123 Smith, E., Lynch, R. S., & Pisano, D. J. 2022, *AJ*,
1124 164, 123, doi: [10.3847/1538-3881/ac7e47](https://doi.org/10.3847/1538-3881/ac7e47)
1125 Thompson, A., & Nicely, M. 2021, *cuSignal: The*
1126 *GPU-Accelerated Signal Processing Library*,
1127 0.19.0. <https://github.com/rapidsai/cusignal>
1128 Vafaei Sadr, A., Bassett, B. A., Oozeer, N.,
1129 Fantaye, Y., & Finlay, C. 2020, *MNRAS*, 499,
1130 379, doi: [10.1093/mnras/staa2724](https://doi.org/10.1093/mnras/staa2724)
1131 Yuan, M., Zhu, W., Zhang, H., et al. 2022,
1132 *MNRAS*, 513, 4787, doi: [10.1093/mnras/stac963](https://doi.org/10.1093/mnras/stac963)