# Radio Astronomical Monitoring in Virtual Environment

Kira Konich[1], Igor Nikitin[2], Stanislav Klimenko[3],
Valery Malofeev[4], and Sergey Tyul'bashev[4]

[1] Bauhaus University, Weimar, Germany
[2] Fraunhofer Institute for Algorithms and Scientific Computing, Sankt Augustin, Germany
[3] Institute of Computing for Physics and Technology, Protvino, Russia
[4] Pushchino Radio Astronomy Observatory, Lebedev Physical Institute, Pushchino, Russia
kira.konycheva@uni-weimar.de, igor.nikitin@scai.fraunhofer.de,
stanislav.klimenko@gmail.com, malofeev@prao.ru, serg@prao.ru

**Abstract**

We present *StarWatch*, our application for real-time analysis of radio astronomical data in Virtual Environment. Serving as an interface to radio astronomical databases or being applied to live data from the radio telescopes, the application supports various data filters measuring signal-to-noise ratio (SNR), Doppler's drift, degree of signal localization on celestial sphere and other useful tools for signal extraction and classification. Originally designed for the database of narrow band signals from SETI Institute (setilive.org), the application has been recently extended for the detection of wide band periodic signals, necessary for the search of pulsars. We will also address the detection of week signals possessing arbitrary waveforms and present several data filters suitable for this purpose.

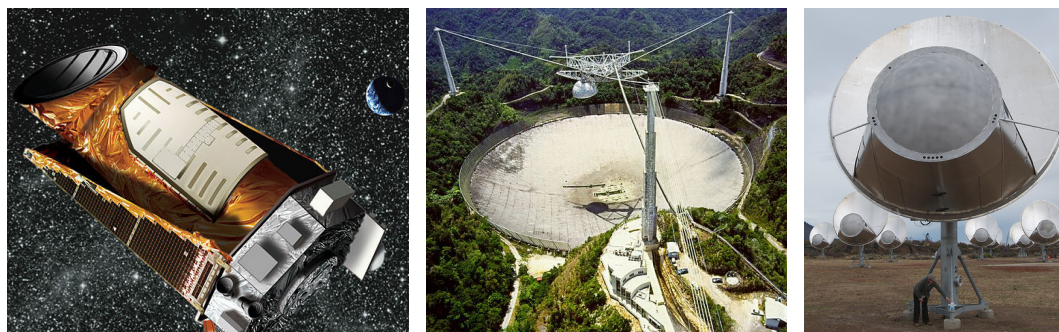*Keywords:* scientific visualization, virtual environments, Big Data, radio astronomy



Figure 1: Is there anybody out there? On the left: telescope Kepler, in the middle: telescope Arecibo, on the right: telescope SETI ATA-42 (courtesy of Colby Gutierrez-Kraybill, www.flickr.com/photos/cgk/1558787110, creativecommons.org/licenses/by/2.0).
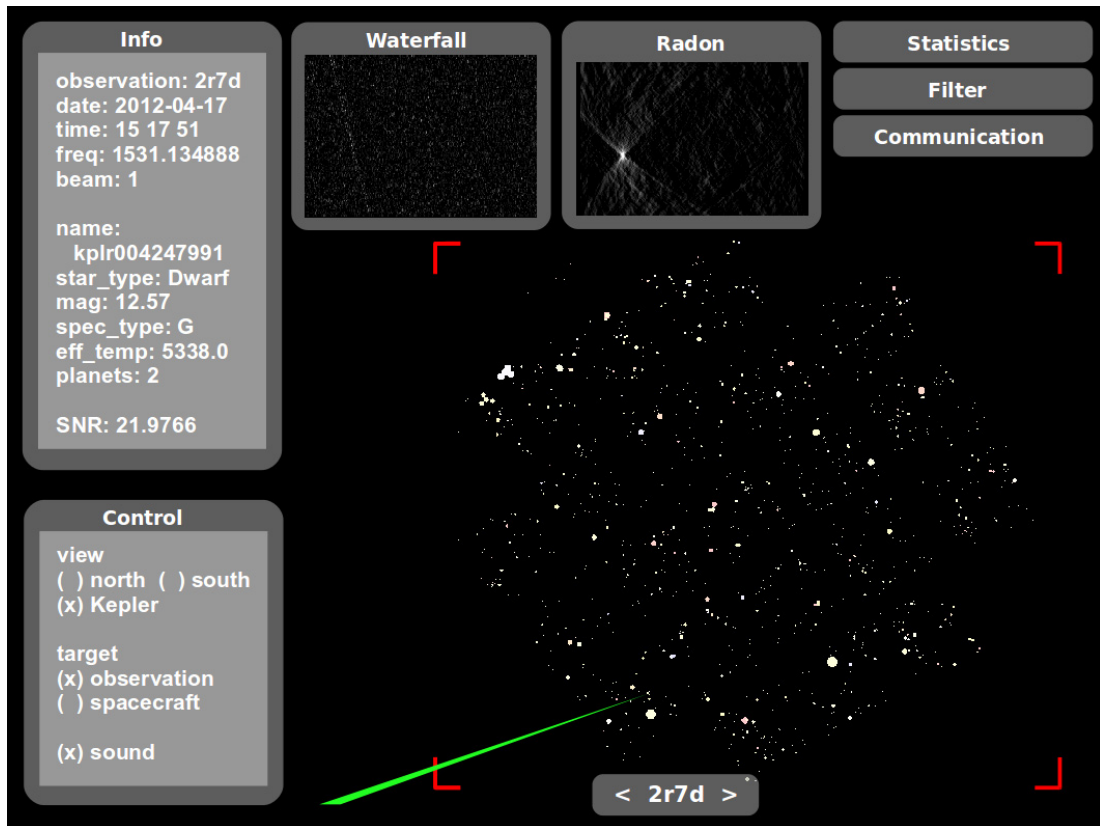
Figure 2: Screenshot of the application *StarWatch* for radio astronomical monitoring in Virtual Environment. The observations are shown as a sky map, here representing a combined survey of the optical telescope Kepler and the radio telescope ATA-42. For a selected observation the corresponding meta-information is displayed and a pointing ray is shown. Color of the ray (here green) indicates that the signal possesses large signal-to-noise ratio, strong localization on celestial sphere and considerable Doppler's drift, i.e. satisfies SETI criteria for potential signals of extraterrestrial origin. Various data filters (Radon transform etc) allow to study the signal in detail.

# 1   Introduction

SETI stands for a Search for Extra-Terrestrial Intelligence, a long-term project conducted by the institute of the same name [1]. The data are acquired from Allen Telescope Array (ATA-42), consisting of 42 antennas in LNSD configuration (Large Number of Small Dishes). The array supports simultaneous data acquisition from several directions on celestial sphere at angular resolution 245"x118" and frequency range 0.5-11.2 GHz.

Sensitivity of the telescopes used in SETI project can be demonstrated using the following estimations. Arecibo radio telescope is the world's largest dish telescope, which can also work in transmission mode. It has 305m diameter and 1MW emitter power. Estimation [2] shows, that emitter with such parameters can be registered by ATA-42 at a distance 150 light years, provided that the emission axis is directed precisely towards the Earth. Estimation [3] for Green

Bank Telescope (GBT), the other telescope used in SETI project, extends this distance to 1000 light years.

Till now the search of extraterrestrial signals produced no results other than a single observation on 15-Aug-1977, so called Wow!-signal [4]. The signal has been detected by a Big Ear, the telescope used in SETI project that time. The signal was located in Sagittarius constellation, close to the hydrogen line frequency, possessed very strong SNR $\sim 30$, lasted 72 seconds and never repeated again.

The main obstacle for detecting the extraterrestrial signals is a problem of large search. The frequency range 1-10GHz forms so called "microwave window" where interstellar transmissions are principally possible. Granularity of the search should be about 1Hz for detection of narrow band signals. This brings a factor $10^{10}$, while the scanning of celestial sphere with 245"x118" tiles brings a factor of $10^7$. Measuring and processing of so large datasets is hard to accomplish. Therefore SETI usually narrows the search to a frequency range 1.4-1.7GHz ("water hole") and restricts the scanning to the positions of extra-Solar planets (exoplanets) marked by optical telescope Kepler. The extension of ATA telescope to 350 dishes is planned. Also SETI affiliates several web projects where the enthusiasts can help to process the large data volumes (setilive.org) and to develop new algorithms of data analysis (setiquest.org).

Setilive.org is a web project forwarding SETI data for the analysis of volunteers. Till 12-Oct-2014 it supported live feeds of signals from ATA-42, which then have been discontinued. Now setilive.org serves as a large archive of radio astronomical data with more than 1.5 millions observations for more than 7.5 thousands observation targets. In this paper we present an interface to these data, which uses interactive Virtual Environments for representation of radio astronomical observations and supports special tools for signal extraction and classification. The methods of signal processing used in this software have been described in our previous papers [5, 6]. Since the input data usually come in form of plain images, the software can be applied to live feeds or other signal collections possessing similar data organization. Typically, the telescope performs 90sec observation, the result is placed in Internet as 80KB PNG image, its processing requires 0.4sec and finally the data are visualized in our application at 60fps rate. In this paper we will extend the methods from the narrow band SETI signals to wide band periodical signals from pulsars and other waveforms.

In Section 2 we overview our application, its main features and the user interface. In Section 3 we present the principal scheme of the application, data structure and its communication with the real world. In Section 4 we go deeper into implementation details and present the structure of data filters used for signal analysis. Section 5 presents the results obtained with the application.

## 2   User interface

The application can be used as full stereoscopic 3D Virtual Environment or in 2D desktop mode, dependently on selected configuration. The distribution of observations in the database is represented as a sky map. Brightness and color are assigned according to the magnitude and spectral class of stars (when available). There are several views: north, south and dedicated ones, those can occupy a smaller region of the sky and represent a special subset of the observations. E.g. SETI uses the positions of exoplanets detected by the orbital telescope Kepler for targeting the Earth-bound telescope ATA-42. Kepler's field of view has a diameter $\sim 12°$ and contains $\sim 4700$ observations of exoplanets (counting confirmed planets and unconfirmed candidates). On the corresponding view Fig.2 one can see the Kepler's observations and even recognize a structure of CCD matrix used in the telescope. In addition to Kepler's exoplanets,

setilive data contain other objects: a number of stars selected from the astronomical catalogues as potentially interesting for SETI project [7, 8], observations of several human-made spacecrafts used for system calibration and some exotic targets like position of Wow!-signal.

The application also contains several panels, which can be opened, closed and moved on necessity. Control panel allows to select the view, targets of observation, switch on/off audio representation of the signal. Info panel contains meta-information on the given observation, such as observation Id, date and time of registration, frequency range, component Id (so called beam number), catalogue name, stellar magnitude, spectral type, temperature and the number of known planets when available. It can also display the signal-to-noise ratio, a schematic view for the spacecrafts and other information associated with the observation.

Waterfall plots show setilive data in their original form, as a frequency spectrum (horizontal axis) varying in time (vertical axis). SETI is mainly interested by narrow-band signals, arguing that concentration of signal power in a narrow <1Hz band will strongly distinguish it from the natural sources (stars, pulsars, quasars, ...) which all have much wider bandwidth. Therefore, if a hypothetical sender of the interstellar message wants it to be received and decoded, it should be sent in a narrow band.

On the waterfall plots the narrow band signals have a form of straight lines possessing a slope to the vertical axis due to Doppler's drift. The drift indicates an accelerated motion, such as motion of the source on a circular orbit. More precisely, in our case it indicates a relative accelerated motion of the source and the Earth-bound telescope. Terrestrial radio sources (TV-stations, mobile phone towers, etc) as well as the satellites on geostationary orbit possess no Doppler's drift and can be excluded from consideration using this criterion. Further, setilive observations come in groups, consisting of simultaneous signal measurements in several near locations on celestial sphere, so called beams. The signals from near-Earth radio sources normally give contribution in all beams. The signals present only in single beam are typically located outside of the Solar System. Therefore, according to SETI argumentation [9], the signal has potentially extraterrestrial origination, if it is (1) concentrated in a narrow band, (2) possesses considerable Doppler's drift and (3) appears as a single beam in a group of simultaneous observations.

In our *StarWatch* application the ray pointing towards the signal source shows classification of the signal to the following three categories: green ray - signal with extraterrestrial signature (narrow band + Doppler's drift + only one beam has large SNR); blue ray - noise (all beams have small SNR); red ray - interference (all other signals).

# 3   Internal structure

The application is created using Avango VE framework (avango.com). Avango in its current version (NG - New Generation) uses Open Scene Graph as graphics engine, Python as scripting layer and C++ as API. Avango application is constructed from precompiled building blocks using connections between data fields. The structure of the application is set via scripting language, while C++ objects support the required processing speed. This approach allows rapid development of utterly complex applications, while the task of 3D visualization is completely overtaken by the framework. To present the application in 3D, one just need to switch the computer to a modern 3D beamer, e.g. the one with DLP-linked shutter glasses. This technique does not require separate IR-transmitters or special screens. Just use a white wall for the projection, switch the application in stereoscopic mode and enjoy 3D immersive visualization.

Fig.3 shows the principal scheme of our application. The source database is typically a web server from which the data can be received via URL queries. *Control* panel is inherited
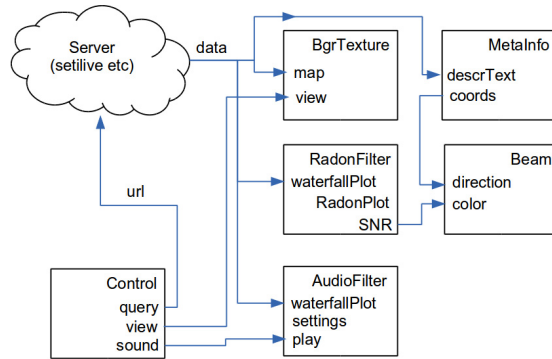
Figure 3: Principal scheme of the application, with activated Radon filter.

from Avango *Panel* base class and supplied with the necessary fields. *BgrTexture* implements the background sky map and switching between the views. *MetaInfo* is one more panel, displaying the text and images associated with a given observation. *Beam* is a ray primitive whose orientation is defined by sky coordinates of the observation and color coding comes from SNR-classification of the beams. Audio-representation of signals is optionally enabled. For this purpose inverse FFT of the waterfall plot is evaluated, the obtained signal is converted to WAV format and played with a suitable sound player. Audio allows to represent such signal parameters as bandwidth, Doppler's drift, modulation etc. Other filters will be now explained in details.

# 4    Data filters

*Radon transform* is used for detection of straight lines on the waterfall plots. It is the integral of the form:
$$R(x, \alpha) = \int dy \, w(x \cos \alpha - y \sin \alpha, x \sin \alpha + y \cos \alpha)$$

where $w(x, y)$ is the waterfall plot and $\alpha$ is the slope parameter. This integral accumulates lines into points and amplifies SNR by a factor $n^{1/2}$ for the images of size $n \times n$. On the corresponding Radon plot Fig.2 the signal has a form of bright spot. After the Radon transform, we measure SNR for all beams in the observation group and use them for signal classification. The numerical evaluation of the integral in Radon transform requires $O(n^3)$ floating point operations. Waterfall plots in SETI observations are typically 768x384 PNG images, rescaled to 256x256 for the purpose of our analysis. Processing of one observation in this setup requires 0.4sec on 3GHz Intel i7 processor. Every such observation represents 1.5min scan of 533MHz frequency band in a certain direction on celestial sphere. This is just a little tile in the global search pattern and the overall effort is multiplied to the number of such tiles.

*Folding transform* is used for detection of generic periodic signals. It has a form:

$$F(\phi, T) = \int dt \, s(t) \Delta(t - \phi, T)$$

where $s(t)$ is a signal, $\Delta(t, T)$ is $T$-periodic Dirac's function and $\phi$ is the phase parameter. The integral accumulates the data with equal phases from various periods and amplifies SNR for

the periodic signals. The amplification factor is $(np/T)^{1/2}$, the granularity of the period $T$ is $dT = T/np$, where $np$ is the number of data points taken into analysis. For smaller $dT$ the algorithm will acquire the same data points and lead to identical results. Here and further $T$, $dT$ and other time related values are measured in data points. The folding transform requires $O(np^2 \ log(T_1/T_0))$ operations, for the period scanned in the range $[T_0, T_1]$.

The folding transformation is more suitable for the signals which are periodic but not localized in a narrow band. In particular, pulsar signals have such a form [14]. They are not harmonic and in frequency spectrum their power is spread in a number of subharmonics. Weak signals of these form can be easily lost in frequency spectrum but can be recovered by folding transform. In more detail, the pulsar signals studied on BSA telescope at Pushchino Radio Astronomy Observatory, are registered simultaneously in a number of beams, typically $nb = 48$, and in a number of frequency bands, typically $nf = 6$ plus one cumulative. The sampling rate is about 10Hz, so that one minute of measurement contains $np = 600$ points and requires about 0.8MB of disk space. There are also "high density" data with more severe characteristics.

Pulsar signals are not sensitive to Doppler's drift, since this effect is much smaller than their bandwidth. There is another effect, dispersion on interstellar medium [15], which leads to the phase shifts between frequency bands and produces the slopes very similar to the Doppler's drift. To amplify SNR, one needs to to compensate these shifts with a *matching algorithm*, a version of Radon transform restricted to few discrete bands. Computational complexity of the matching algorithm is $O(np \ dpmax \ (T_1 - T_0))$, where $dpmax$ is a maximal dispersion shift, in points per frequency band. In the setup above typically $dpmax = 14$. This effort is additive to the folding transform, the both should be also multiplied to the number of beams and frequency bands.

One more specifics about BSA data is that the telescope beams are continuously sweeping across the sky and the pulsar signals appear in data during a restricted time interval, when the beam is directed precisely on the pulsar. Typically these intervals have 1-5min duration and the search should be done with the overlapping, e.g. by taking 5min segment and sequentially shifting it for 1min along time axis. The overlapping additionally increases the computational effort, so that processing of BSA data in average requires 18sec per a minute of measurements.

Mass processing of SETI and pulsar data can be trivially parallelized to make use of all available cores and processors, so that the processing speed can be significantly increased on parallel architectures.

Now we see that the choice of the data filter depends on the form of the signal: narrow band SETI signals can be processed by Radon transform, wide band pulsar data are better to process with folding transform. A question remains what to do with the unknown waveforms. What if senders from some of their reasons will use not a narrow band or a pulse mode of communication but a signal of entirely other form? The example can be so called noise-like signal [10], which cannot be distinguished from the noise by its spectrum. In constructive form the question can be reformulated as follows:

*Is there an algorithm to distinguish a given time series from a random sequence?*

Here we present the ideas for two such algorithms.

*Entropy of the signal.* It is well known that the entropy is a measure of chaos. Being taken with the opposite sign, negentropy is a measure of information content. In more pragmatic view the entropy can be considered as a measure of random uniformity of a sequence. In continuous case

$$E[\rho] = \int d^n x \ f(\rho(x)), \quad f(\rho) = -\rho \log \rho$$

where $E$ is the entropy, $\rho(x)$ is the probability density, $x$ is $n$-dimensional vector parameter. Taking a variation of this functional on the surface of normalization constraint $\int dx\,\rho = 1$, one can show that this functional achieves maximum on a uniform distribution $\rho = Const$. Thus, the entropy can be used to measure a deviation of the probability density from the uniform distribution. Actually, this holds for any such functional with $f'' < 0$, not only for the entropy.

For time series (value vs time) one usually interests not by the entropy of the value itself. E.g. in language analysis different letters possess different frequency of appearance. The probability distribution of letters is not uniform and the entropy of letter is not maximal. However this says nothing about the informational content of the message. It makes sense to consider $n$ sequential letters as one word and to measure the entropy in the corresponding $n$-dimensional space of words. In this way one can distinguish a meaningful text from a random permutation of its letters. Also in our case, not the instant value of the signal is important, but the way of its variation in time, which can encode different waveforms (harmonic, pulse, others).

The eventually non-uniform distribution of a single value can be easily corrected by the following *flattening algorithm*. One constructs a cumulative distribution function (CDF) for a sequence by its sorting in ascending order. Further, the actual values in the original sequence are replaced with their numbers in the sorted sequence, normalized to the total length of the sequence. The obtained sequence of real numbers is now uniformly distributed in the range $[0, 1]$. It can be used in further analysis on the place of the original sequence, since it has the same ordering of the values and the same shape of the waveforms. If one does not interest in fine details of the waveforms, one can quantize the interval by $q$ equidistant levels and round continuous values to these discrete levels.

Direct computation of multidimensional integral in entropy definition suffers from the problem known as *curse of dimensionality*. Assuming $q$ quantized intervals, in $n$-dimensional space one has $q^n$ cubes. Every cube should be populated by many points to have a stable definition of probability density. To achieve this, one must have a large sequence length $N >> q^n$. Practically, one can use only small $q$ and $n$, e.g. for $N = 10^{10}$ use $q = 10$ and $n = 8$.

*Cryptographic strength.* The idea to subject cosmic noise to cryptographic analysis sounds really funny. After all, there is nothing to decrypt there! However, there are cryptographic tests for random number generators which measure no more no less then the randomness of a given numerical sequence. Therefore, to test whether a measured signal resembles a random sequence, it's sufficient to pass it through a standard test for random number generators.

More specifically, a bit sequence is called cryptographically strong, if prediction of the next bit from the previous bits of the sequence with any polynomial-time algorithm has success ratio not better than a wild guess (50%). The test is archetypic in the sense that all other polynomial-time statistical tests for randomness can be reduced to the next-bit test in polynomial time [11]. In practice, there are tests for random number generators known as *DieHard* [12] and *DieHarder* [12]. They include a collection of subtests returning so called p-values, random variables uniformly distributed in the range $[0, 1]$ whenever the input sequence is uniformly random. Further one performs *Kolmogorov-Smirnov test* for the uniformity of the distribution of p-values, by forming CDF of p-values and measuring its maximal deviation $D_p$ from the linear function. Kolmogorov's distribution can be used to estimate the confidence levels, e.g. $99.73\%CL$ corresponds to $D_p < 1.81 n_p^{-1/2}$, where $n_p$ is a number of p-values.

To compare the efficiency of various filters, we have prepared three synthetic data samples: narrow band signal of SETI type, periodical pulsar spikes and free waveform signal, for which we use WAV sample of human speech. The clean samples are linearly mixed with a random noise: $s = a$ sample $+ (1 - a)$ rnd, where $a \in [0, 1]$, so that input SNR $= a/(1 - a)$. The

obtained signals are processed by the data filters. Fig.4 shows the response of the filters as a function of input SNR. For Radon and folding filters the sample size $N = 10^4$ is taken. Both filters start to work at SNR > 0.1, Radon a bit better for narrow band SETI signals, folding with equal efficiency for SETI signals and pulsars. Speech waveforms are recognized at SNR > 1 by both filters. Similar results are produced by the entropy filter, for which a larger sample $N = 10^5$ is taken and the internal parameters are set to $q = 2$, $n = 8$. The figure presents information density $I = 8 - E/log(2)$, in bits per symbol. In this plot small SNR correspond to maximal entropy 8bits/symbol and zero information content, while increasing SNR indicate larger information content up to 0.3-6bits/symbol. Cryptographic DieHard filter requires at least $N = 10^8$ sequence to begin with. So big datasets are already unavailable in the above mentioned sources, while the project setiquest.org has raw data of even larger size. In our synthetic examples we generate the waveforms indefinitely, in particular repeat 10sec speech sequence $10^4$ times on non-periodic random background. The plot shows the results of Kolmogorov-Smirnov test with marked 99.73% confidence level. Data below this level can be considered as random noise, while above this level the filter indicates the presence of a signal.
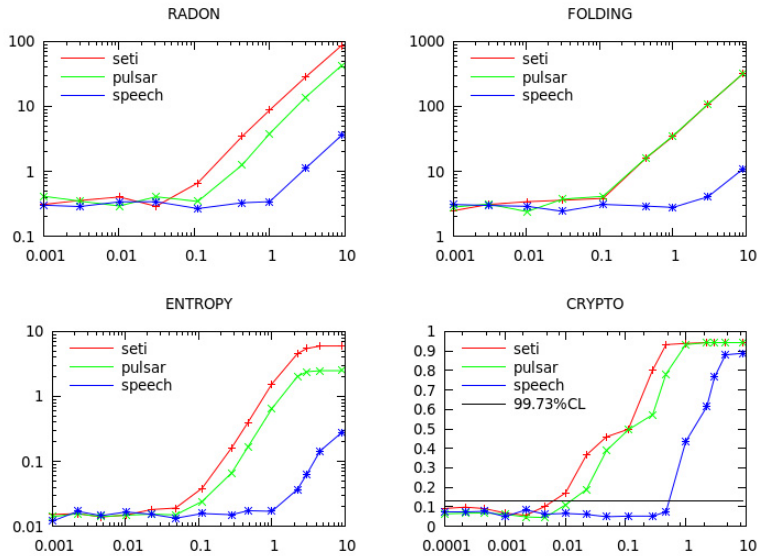


Figure 4: Response of data filters as a function of input SNR.

# 5   The results

In our previous work [6] we have performed statistical analysis of setilive data, using Radon transform and specially constructed filter for selection of single beams. With these means we have found 28 strong single beam signals. Some of these signals are shown on Fig.5. In addition, we have performed statistical estimation of signal background using Monte Carlo technique and marked 1072 statistically significant single beam signals visible above the background.

The selected signals should be subjected to more sophisticated analysis [3]. Although all of them are concentrated in a narrow band, possess Doppler's drift and appear only in one beam, there is still a small probability that a satellite signal comes to the telescope through so called sidelobe [9] and will be recorded as a single beam event. To exclude this possibility, a
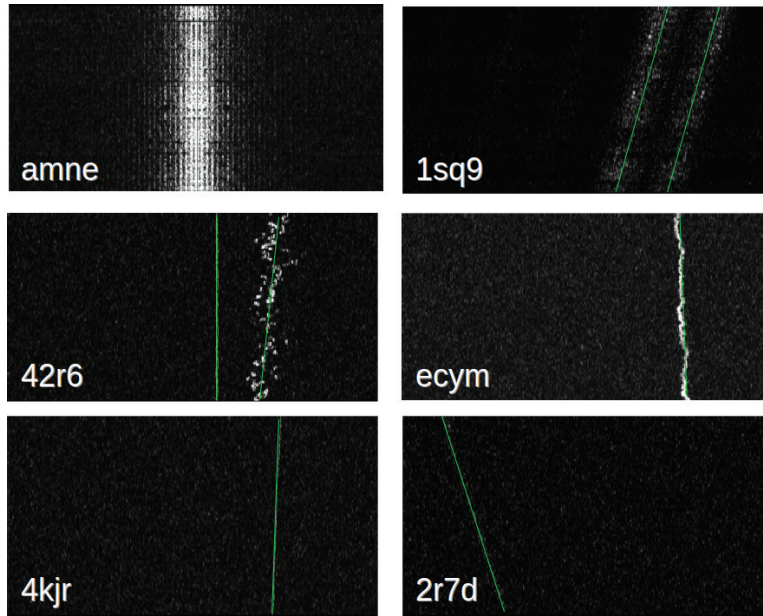
Figure 5: Typical setilive signals.   The figure shows  waterfall plots – frequency spectrum (horizontal axis) varying in time (vertical axis): (amne) terrestrial source; (1sq9) GPS satellites; (42r6) combination of terrestrial and satellite signals; (ecym) satellite signal, wavy form indicates own rotation; (4kjr) single beam - potential signal of ET origination; (2r7d) one more single beam. The plots can be retrieved by their 4-character code from <talk.setilive.org/observation_groups/GSL000****>.
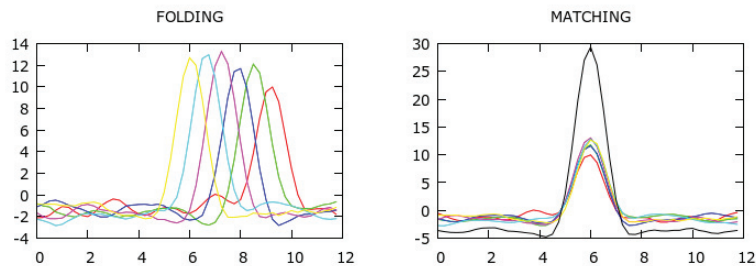


Figure 6: Typical pulsar signal from BSA telescope. On the left: the result of folding transform, the peaks in 6 frequency bands arrive at different time due to dispersion on interstellar medium; on the right: matching transform collects the peaks together and increases total SNR. The figure shows SNR as a function of time, in data points.

cross-validation procedure should be performed: registration of signals with similar parameters (frequency, bandwidth, modulation) in different directions would indicate that it is a signal from a satellite which must be rejected. Such tests require intensive comparisons across the database and we are planning them in our future research.

The analysis of pulsar data from Pushchino Radio Astronomy Observatory is now running

in test phase. With our implementation of folding and matching algorithms the data from 1 hour measurement have been processed. All pulsars known to be present in the data at the level SNR $> 5$ have been extracted. Also one pulsar at the level SNR $\sim 3$ has been found, which has not been detected in these data by previously used methods. The typical image of extracted pulsar signal is shown on Fig.6.

We plan to improve the algorithms using fast Radon transform [16], fast folding [17], other accelerated techniques [18] and apply them to high density pulsar data.

In this paper we have also performed a proof of concept for two generic waveform algorithms. As expected, they require a bigger data samples than the algorithms for fixed waveforms. On the other hand, in synthetic examples they have shown the ability to detect the signals of unknown form. Further we plan to test these algorithms on real radio astronomical data.

# References

[1] Don Osgood, Ronald D. Ekers, SETI 2020: A Roadmap for the Search for Extraterrestrial Intelligence, SETI Press 2002 .

[2] Andrew P. V. Siemion et al, Searching for Extraterrestrial Intelligence with the Square Kilometre Array, <arxiv.org/abs/1412.4867>

[3] Andrew P. V. Siemion et al, A 1.1 to 1.9 GHz SETI Survey of the Kepler Field: I. A Search for Narrow-band Emission from Select Targets, The Astrophysical Journal, 767:94 (13pp.), 2013, <arxiv.org/abs/1302.0845>

[4] Jerry R. Ehman, The Big Ear Wow!-Signal: What We Know and Don't Know About It After 20 Years, <www.bigear.org/wow20th.htm>

[5] S.V.Klimenko et al, On signals with Doppler's drift, fast Fourier transform and search for extraterrestrial intelligence, in Proc. of SCVRT2013, Protvino, Russia, November 26-28 2013.

[6] S.V.Klimenko, I.N.Nikitin, On statistical data accumulation, Radon transform and search for extraterrestrial intelligence, in Proc. of CPT2014, Cyprus, Larnaca, May 11-18, 2014.

[7] Margaret C. Turnbull, Jill C. Tarter, Target Selection for SETI: 1. A Catalog of Nearby Habitable Stellar Systems, The Astrophysical Journal Supplement Series, 145:181198, 2003.

[8] Margaret C. Turnbull and Jill C. Tarter, Target Selection for SETI. II. Tycho-2 Dwarfs, Old Open Clusters, and the Nearest 100 Stars, The Astrophysical Journal Supplement Series, 149:423436, 2003.

[9] SETI Tutorials: Signal Processing and SETI <setiquest.org/about/tutorials>, The Doppler Effect <setilive.org/about>, How would we know that the signal is from ET? <www.seti.org/faq>

[10] Don Torrieri, Principles of Spread-Spectrum Communication Systems, Springer 2015.

[11] Andrew Chi-Chih Yao. Theory and applications of trapdoor functions, In Proc. of the 23rd IEEE Symposium on Foundations of Computer Science, 1982.

[12] George Marsaglia, The Marsaglia Random Number CDROM including the Diehard Battery of Tests of Randomness, <www.stat.fsu.edu/pub/diehard>

[13] Robert G. Brown, Dirk Eddelbuettel, David Bauer, Dieharder: A Random Number Test Suite Version 3.31.1, <www.phy.duke.edu/~rgb/General/dieharder.php>

[14] Malofeev V.M., Malov O.I., Detection of Geminga as a radiopulsar, Nature V389(1997) pp.697-699.

[15] Kardashev N.S. et al, Review of scientific topics for the Millimetron space observatory, Physics-Uspekhi 57 (2014), pp.1199-1229.

[16] Chandra S., Finite Transform Library (FTL), <finitetransform.sourceforge.net>

[17] David H. Staelin, Fast Folding Algorithm for Detection of Periodic Pulse Trains, Proceedings of the IEEE, 57 (1969).

[18] D. Kent Cullers, Stanley R. Deans (Eds.), SETI Algorithms, Chap.4, The DADD (Doubling Accumulation Drift Detection) algorithm, <ftp://ftp.seti.org/gharp/SetiAlgorithms-CullersDeans.pdf>