



Keitel, D. and Ashton, G. (2018) Faster search for long gravitational-wave transients: GPU implementation of the transient F-statistic. *Classical and Quantum Gravity*, 35(20), 205003. (doi:[10.1088/1361-6382/aade34](https://doi.org/10.1088/1361-6382/aade34))

There may be differences between this version and the published version. You are advised to consult the publisher's version if you wish to cite from it.

<http://eprints.gla.ac.uk/168905/>

Deposited on: 17 September 2018

Enlighten – Research publications by members of the University of  
Glasgow  
<http://eprints.gla.ac.uk>

ACCEPTED MANUSCRIPT

## Faster search for long gravitational-wave transients: GPU implementation of the transient F-statistic

To cite this article before publication: David Keitel *et al* 2018 *Class. Quantum Grav.* in press <https://doi.org/10.1088/1361-6382/aade34>

### Manuscript version: Accepted Manuscript

Accepted Manuscript is “the version of the article accepted for publication including all changes made as a result of the peer review process, and which may also include the addition to the article by IOP Publishing of a header, an article ID, a cover sheet and/or an ‘Accepted Manuscript’ watermark, but excluding any other editing, typesetting or other changes made by IOP Publishing and/or its licensors”

This Accepted Manuscript is © 2018 IOP Publishing Ltd.

During the embargo period (the 12 month period from the publication of the Version of Record of this article), the Accepted Manuscript is fully protected by copyright and cannot be reused or reposted elsewhere.

As the Version of Record of this article is going to be / has been published on a subscription basis, this Accepted Manuscript is available for reuse under a CC BY-NC-ND 3.0 licence after the 12 month embargo period.

After the embargo period, everyone is permitted to use copy and redistribute this article for non-commercial purposes only, provided that they adhere to all the terms of the licence <https://creativecommons.org/licenses/by-nc-nd/3.0>

Although reasonable endeavours have been taken to obtain all necessary permissions from third parties to include their copyrighted content within this article, their full citation and copyright line may not be present in this Accepted Manuscript version. Before using any content from this article, please refer to the Version of Record on IOPscience once published for full citation and copyright details, as permissions will likely be required. All third party content is fully copyright protected, unless specifically stated otherwise in the figure caption in the Version of Record.

View the [article online](#) for updates and enhancements.

# Faster search for long gravitational-wave transients: GPU implementation of the transient $\mathcal{F}$ -statistic

David Keitel<sup>1</sup> and Gregory Ashton<sup>2,3</sup>

<sup>1</sup>University of Glasgow, School of Physics and Astronomy, Kelvin Building,  
Glasgow G12 8QQ, Scotland, United Kingdom

<sup>2</sup>Max Planck Institut für Gravitationsphysik (Albert Einstein Institut), 30161  
Hannover, Germany

<sup>3</sup>Monash Centre for Astrophysics, School of Physics and Astronomy, Monash  
University, VIC 3800, Australia

E-mail: david.keitel@ligo.org

LIGO-P1800031-v6 [draft version: 27 August 2018]

**Abstract.** The  $\mathcal{F}$ -statistic is an established method to search for continuous gravitational waves from spinning neutron stars. Prix et al. [1, (2011)] introduced a variant for transient, hours–months long, quasi-monochromatic signals. Possible astrophysical scenarios for such transients include glitching pulsars, newborn neutron stars and accreting systems. Here we present a new implementation of the transient  $\mathcal{F}$ -statistic, using pyCUDA to leverage the power of modern graphics processing units (GPUs). The obtained speedup allows efficient searches over much wider parameter spaces, especially when using more realistic transient signal models including time-varying (e.g. exponentially decaying) amplitudes. Hence, it can enable comprehensive coverage of glitches in known nearby pulsars, improve the follow-up of outliers from continuous-wave searches, and might be an important ingredient for future blind all-sky searches for unknown neutron stars.

## 1. Introduction

Spinning neutron stars (NSs), when non-axisymmetrically deformed, emit weak but potentially detectable gravitational waves (GWs) [2]. Many searches [3] with the LIGO and Virgo detectors [4, 5] focus on continuous wave (CW) signals that are persistent over a whole observation run, but there are also scenarios for shorter signals from transiently perturbed NSs. If those signals are slowly evolving in frequency and last on the time scale of hours to months, analysis methods adapted from CW searches are well suited to their detection. In [1] (hereafter also referred to as ‘PGM’), the astrophysical motivation for such transient signals was discussed and a matched-filter search method proposed. It is based on the established  $\mathcal{F}$ -statistic, which was introduced in [6, 7] and used in many CW searches [recently e.g. in 8–10].

Matched-filter searches for weak signals from unknown sources (or those with imperfectly known parameters) are computationally very expensive since a wide parameter space needs to be densely covered with templates. Starting from a typical CW search that covers a certain parameter space in signal frequency, spindown and sky location but assumes a constant signal amplitude, the addition of new unknown parameters to describe the transient evolution further increases computational cost.

However, the attractiveness of the transient  $\mathcal{F}$ -statistic algorithm from [1] is that it starts from time-discretised quantities already computed for the standard CW  $\mathcal{F}$ -statistic and then only needs to take partial sums of these to study the set of possible transient signals. Still, for long total observation times the evaluation of these partial sums can easily dominate over the original computational cost, especially if the templates have a non-trivial amplitude evolution. For example, exponential amplitude decay might be expected from a transiently excited and then viscously damped oscillation mode, or for GW emission associated with the observed exponential relaxation phases after pulsar glitches [11, 12].

The task of multiple partial sums of some input data can obviously benefit from massive parallelisation. Here we present a straightforward translation of the algorithm from [1] to run in parallel on the many subprocessors of modern graphics processing units (GPUs), which are addressed through the general-purpose open-source GPU computing package `pyCUDA` [13]. The full data analysis is implemented within the open-source `PyFstat` package [14, 15] for  $\mathcal{F}$ -statistic-based GW searches.†

In the following, we briefly review the formalism from [1] to define the transient  $\mathcal{F}$ -statistic (section 2), then describe its `pyCUDA` implementation (section 3). We test the speed and memory requirements (section 4) and compare with the original CPU implementation from `LALSuite` [16]. The paper ends with a brief discussion (section 5) of how the achieved speedup widens the scope of feasible searches for long CW-like gravitational wave transients. This includes enabling a comprehensive coverage of glitch events in nearby known pulsars, improving the sensitivity of all-sky CW searches through following up more outliers with transient analyses, and the potential use as an ingredient in future blind all-sky searches for unknown disturbed NSs.

## 2. Formalism

We present a straightforward `pyCUDA` implementation of the PGM ‘atoms-based’ transient  $\mathcal{F}$ -statistic algorithm from Appendix A1 of [1]. It is based on a discretised method to compute the overall (CW)  $\mathcal{F}$ -statistic, introduced in [17] and described in detail in [18].

The  $\mathcal{F}$ -statistic (for transient or continuous signals) is essentially a likelihood-ratio test for a time series  $x(t)$ , comparing a signal hypothesis

$$\mathcal{H}_{\text{tS}} : x(t) = n(t) + h(t, \lambda, \mathcal{A}, \mathcal{T}) \quad (1)$$

against the alternative hypothesis of pure Gaussian noise,

$$\mathcal{H}_{\text{G}} : x(t) = n(t). \quad (2)$$

The waveform model  $h(t, \lambda, \mathcal{A}, \mathcal{T}) = \varpi(t, \mathcal{T}) h(t, \lambda, \mathcal{A})$  for a slowly-evolving signal separates into a transient window function  $\varpi(t, \mathcal{T})$  [first introduced in 1] and the standard CW waveform  $h(t, \lambda, \mathcal{A})$  [6, 18]. The latter depends on a set of phase evolution parameters  $\lambda = \{\alpha, \delta, f, \dot{f}, \ddot{f}, \dots\}$  (sky position, frequency, and frequency derivatives or ‘spindowns’) and on four amplitude parameters  $\mathcal{A} = \{h_0, \cos \iota, \psi, \phi_0\}$ . ( $h_0$  is a dimensionless amplitude,  $\iota$  and  $\psi$  describe the orientation and polarisation of the source, and  $\phi_0$  the GW phase at a reference frequency. Also see [1, 6, 18] for details on how the  $\mathcal{F}$ -statistic analytically maximises over these parameters.) For the transient part, we currently consider either rectangular or exponential window functions, with

† Latest `PyFstat` source code and examples also available from:  
<https://gitlab.aei.uni-hannover.de/GregAshton/PyFstat/>.

GPU transient  $\mathcal{F}$ -statistic

3

parameters  $\mathcal{T} = \{t_0, \tau\}$  where  $t_0$  is the start time of a signal and  $\tau$  is a duration parameter:

$$\varpi_{\text{rect}}(t, t_0, \tau) := \begin{cases} 1 & \text{if } t \in [t_0, t_0 + \tau] \\ 0 & \text{otherwise,} \end{cases} \quad (3)$$

$$\varpi_{\text{exp}}(t, t_0, \tau) := \begin{cases} e^{-(t-t_0)/\tau} & \text{if } t \in [t_0, t_0 + 3\tau] \\ 0 & \text{otherwise.} \end{cases} \quad (4)$$

The cutoff of  $\varpi_{\text{exp}}$  at  $3\tau$  was introduced in [1] as a speedup optimisation in the understanding that the signal-to-noise ratio (SNR) after this point will be negligible.

The (transient)  $\mathcal{F}$ -statistic is then proportional to the log odds between  $\mathcal{H}_{\text{ts}}$  and  $\mathcal{H}_{\text{G}}$ , after maximising over  $\mathcal{A}$  (or marginalising, see [1, 19, 20] for details):

$$e^{\mathcal{F}(x, \lambda, \mathcal{T})} \propto \frac{P(\mathcal{H}_{\text{ts}}|x, \lambda, \mathcal{T}, \mathcal{I})}{P(\mathcal{H}_{\text{G}}|x, \mathcal{I})}. \quad (5)$$

It can be written as

$$\mathcal{F}(x, \lambda, \mathcal{T}) = \frac{1}{2} x'_\mu(\lambda, \mathcal{T}) \mathcal{M}'^{\mu\nu}(\lambda, \mathcal{T}) x'_\nu(\lambda, \mathcal{T}), \quad (6)$$

where the indices  $\mu, \nu$  run over the four amplitude parameters  $\mathcal{A}$ ,  $\mathcal{M}'^{\mu\nu}$  is the antenna pattern matrix,  $x'_\mu$  are projections of the data onto the model waveforms, and the prime denotes transient windowing. (See Eqs. (32–36) of [1].)

The standard algorithm used in  $\mathcal{F}$ -statistic searches for continuous signals splits a data set starting at  $T_0$  and of length  $T_{\text{obs}}$  into several Short Fourier Transforms (SFTs) of length  $T_{\text{SFT}}$ . [18] describes how to approximate (6) from the per-SFT, per-detector discretised versions of  $\mathcal{M}'^{\mu\nu}$  and  $x'_\mu$ ; in practice we consider the equivalent set of quantities  $\{a_j, b_j, F_{aj}, F_{bj}\}$  as the *atoms* of our  $\mathcal{F}$ -statistic computation, where the  $j$  index runs over SFTs.

The  $a_j$  and  $b_j$  atoms are summed up to yield the discretised antenna pattern matrix elements  $\hat{A}, \hat{B}, \hat{C}$  [defined in Eq. (130) of 18] and their determinant  $\hat{D} = \hat{A}\hat{B} - \hat{C}^2$ , and together with the summed data-dependent complex quantities  $F_a, F_b$  [Eq.(129) of 18] they yield the  $\mathcal{F}$ -statistic as:

$$\mathcal{F}(x, \lambda, \mathcal{T}) = \hat{D}^{-1} (\hat{B} [\Re^2(F_a) + \Im^2(F_a)] + \hat{A} [\Re^2(F_b) + \Im^2(F_b)] - 2\hat{C} [\Re(F_a)\Re(F_b) + \Im(F_a)\Im(F_b)]). \quad (7)$$

(All quantities on the right hand side are understood as depending on  $\lambda$  and  $\mathcal{T}$ , too.)

For persistent CWs, this is evaluated summing all atoms over the full  $T_{\text{obs}}$ . To search for transient signals, we define a grid in  $\{t_0, \tau\}$  space indexed by  $m$  for the  $t_0$  dimension and  $n$  for the  $\tau$  dimension. We indicate the resolutions of this grid as  $dt_0$  and  $d\tau$ ; a natural choice is  $dt_0 = d\tau = T_{\text{SFT}}$  though a coarser or even variable sampling is also possible. Then our goal is to compute, for each  $\lambda$  and a specific window choice  $\varpi$ , the matrix

$$\mathcal{F}_{mn}(\lambda) := \mathcal{F}(x, \lambda, \varpi, t_{0m}, \tau_n), \quad (8)$$

which we also refer to as the *transient  $\mathcal{F}$ -statistic map*. Computing it this way is convenient because the set of atoms  $\{a_j, b_j, F_{aj}, F_{bj}\}$  is only computed once, over the full  $T_{\text{obs}}$ , and this is already done for the CW  $\mathcal{F}$ -statistic anyway. Subsequently, the transient  $\mathcal{F}_{mn}$  map is obtained by evaluating (7) for partial sums of the atoms.

### 3. Implementation

The new GPU version of the transient  $\mathcal{F}$ -statistic is implemented in the framework of the `PyFstat` python package [15]. `PyFstat` has primarily been developed for the follow-up of CW candidates with Markov-chain Monte Carlo (MCMC) methods [14]. Through C-to-python wrappers [SWIG, 21, 22], it uses standard CW search functionality (written in C) from the `LALPulsar` package of the LIGO collaboration's algorithm collection `LALSuite` [16]. Hence, `PyFstat` can also be used as a convenient way to develop modular custom searches for CWs and similar signals, e.g. for the long-duration transients we consider here.

For the transient  $\mathcal{F}$ -statistic, the `PyFstat`-based search application first calls the standard `LALPulsar` algorithm `ComputeFstat` for computing the CW  $\mathcal{F}$ -statistic over the whole data set [17, 18]§. This already takes care of reading the data SFTs, the signal-parameter-dependent translation between detector frames and a common solar system reference frame ('barycentring'), and computation of the per-SFT matched-filter atoms. While for a standard CW search, the atoms would be discarded after computing the overall  $\mathcal{F}$ -statistic, the only change for our transient application is that we ask the `ComputeFstat` routine to also return these atoms for further processing.

The input data for computing the transient  $\mathcal{F}$ -statistic map  $\mathcal{F}_{mn}$  from (8) consists then of only the atoms (a set of vectors of  $N_{\text{SFT}}$  elements each) and the parameters describing a transient window function and grid in  $\{t_0, \tau\}$  space. The atoms, given as 3 real vectors  $a^2(t)$ ,  $b^2(t)$  and  $a(t) \cdot b(t)$  and 2 complex vectors  $F_a(t)$ ,  $F_b(t)$ , are transferred to the GPU as a  $7 \times N_{\text{SFT}}$  real matrix.

The basic idea of massively-parallelised computation on a GPU is to run a grid of identical kernels, each processing the subset of data identified by the kernel's (multi-)index. We provide two structurally different kernels and grid setups for rectangular and exponential windows. To account for the general case where resolutions in  $t_0$  or  $\tau$  different from  $T_{\text{SFT}}$  might be desirable, or where there are gaps in the data, we use  $N_{t_0}$  and  $N_\tau$  for the number of grid points in each dimension, which need not be equal to each other nor to  $N_{\text{SFT}}$ .

In the **rectangular case**, an obvious optimisation was already pointed out in [1] and is implemented in `LALPulsar`: For each starting time  $t_{0m}$ , one can compute  $\mathcal{F}_{mn}$  for all durations  $\tau_n$  by keeping the partial sums of each atom up to each  $\tau_{n'}$  in memory and only adding the atoms with index  $n' + 1$  in the next step. It would thus be wasteful to run a full  $N_{t_0} \times N_\tau$  grid of kernels on the GPU, and instead we only launch  $N_{t_0}$  kernels, each of which internally loops over  $\tau$  and keeps the partial sums in local memory.

In the **exponential case**, no such simple trick is possible, since the contribution to each partial sum at each timestep includes amplitude-weight factors (see Eq. (4)) depending on the  $\tau$  currently being evaluated. Hence, we employ a brute-force grid of  $N_{t_0} \times N_\tau$  kernels on the GPU, each of which only computes the partial sums for a single  $\mathcal{F}_{mn}$ . This grid structure would also be appropriate for any other generic window function.

In both cases, the last steps, still done inside the GPU kernel, are to compute the antenna pattern matrix determinant  $\hat{D}$  and the transient  $\mathcal{F}_{mn}$ -statistic from Eq. (7).

§ As of the writing of this paper, documentation is available at:  
[https://lscsoft.docs.ligo.org/lalsuite/lalpulsar/group\\_\\_\\_compute\\_fstat\\_\\_h.html](https://lscsoft.docs.ligo.org/lalsuite/lalpulsar/group___compute_fstat__h.html)

#### 4. Tests

In this section, we describe tests of the speedup obtained with the `pyCUDA` version, its memory requirements, and its numerical faithfulness to the original implementation. Here we use shorthands ‘rect’ and ‘exp’ for the window functions.

##### 4.1. Speed

We have tested the speed of the `pyCUDA` implementation relative to the standard `LALPulsar` code on several systems. These all have Intel CPUs: a laptop with a Core i5-6200U at 2.30 GHz, a workstation with a Xeon X5675 at 3.07 GHz and two LIGO Caltech cluster nodes with Xeons E5-2630 and E5-2650 at 2.20 GHz each. Note that the `LALPulsar` code is compiled with the aggressive ‘-O3’ optimisation level of the `gcc` compiler for all tests, but it runs on a single core of each CPU only. The `pyCUDA` code was benchmarked on several GPUs from the Nvidia GeForce GTX family (1050, 1060, 1070 and 1080Ti, with 2–11 GB RAM) and on a Nvidia Tesla V100-PCIE (16 GB RAM), all installed on the same workstation and cluster nodes.

We consider observation times  $T_{\text{obs}}$  from 1 hour up to 1 year, with no gaps in the data. Gaussian noise and a transient signal with  $\tau = 0.5 T_{\text{obs}}$  are simulated through `PyFstat`, though the speed of calculating  $\mathcal{F}$ -statistics does not depend on whether the data contains a signal. The SFTs are taken at  $T_{\text{SFT}} = 1800$  s and  $\mathcal{F}_{mn}$  is sampled at  $dt_0 = d\tau = T_{\text{SFT}}$  over a grid of  $t_0 \in [T_0, T_{\text{obs}} - 2T_{\text{SFT}}]$  and  $\tau \in [2T_{\text{SFT}}, T_{\text{obs}}]$ . The upper limit on  $t_0$  and lower limit on  $\tau$  are set because the low-level implementation requires at least 2 SFTs per  $\mathcal{F}_{mn}$  computation.

Since GPU results for a single-template (fixed  $\lambda$ ) analysis might be too pessimistic because of startup overheads, and in practice speedups are only relevant for searches over broad  $\lambda$  regions, we time searches over 100 frequency bins; though for simplicity we assume a fixed sky location and no spindown. Timing results are summarised in figure 1, as average runtime per  $\lambda$  template. As an additional cross-check, these results also include some runs at 1000 frequency bins, which yield consistent timings per template. Note that this is the total runtime of the search (per template), including the initial `LALPulsar` computation of the atoms which always runs on the CPU.

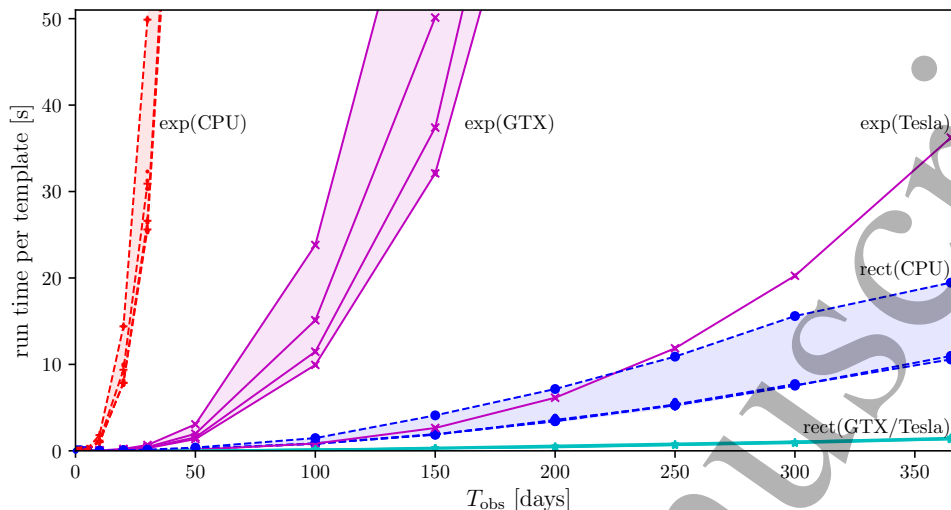
For each architecture, exp windows are much slower than rect windows. As discussed in section 3, for rect windows the number of summation steps can be truncated, while exp windows correspond to the generic case where every step needs to be fully executed. Also, each step for a rect window is a simple addition, while for an exp window evaluation of the actual `exp()` function at each step adds many more floating point operations.

Comparing CPUs with GPUs, we find that the `pyCUDA` version provides speedups of at least an order of magnitude on GPUs of the Geforce GTX 10x0 family compared to the original single-core `LALPulsar` code on contemporaneous CPUs, both for exp and rect windows. In the exponential case, the Tesla V100 provides another similar jump in speed over the GTX family, bringing the cost of exp window transient searches over hundreds of days down to a similar cost as with the standard rect `LALPulsar` CPU implementation.

|| There is no strong incentive to develop a multi-core CPU version, since CW (and transient-CW) search problems can be trivially parallelised by splitting the parameter space into  $N$  sub-regions to run on  $N$  cores, or independent machines.

GPU transient  $\mathcal{F}$ -statistic

6



**Figure 1.** Timing results for both rectangular and exponential transient windows, from CPU (`LALPulsar`) and GPU (`pyCUDA`) implementations on various devices. The vertical axis gives the average run time per template  $\lambda$ . (Most test runs used 100 frequency bins, and a few used 1000 to check the consistency of averages.) Each solid/dashed line connects results from a specific implementation on a specific device, averaging over 3 or more runs at fixed  $T_{\text{obs}}$ , and background shading indicates a specific window run on a family of related architectures. The `exp(CPU)` and `rect(CPU)` families collect results from the four different systems mentioned in Sec. 4.1, while `exp(GTX)` labels results from different Nvidia Geforce GTX 10x0 family devices, the single line labelled `exp(Tesla)` is from a Nvidia Tesla V100, and the `rect(GTX/Tesla)` results are plotted together as they are not significantly different.

We can also more directly compare these measurements to the timing model from Appendix A3 of [1]. We find that to cover arbitrary combinations of  $\{T_{\text{data}}, N_{t_0}, N_\tau\}$ , we need to somewhat generalise the model. This is done in detail in our Appendix A and leads to a timing model with several contributions proportional to the number of data units  $N_{\text{SFT}}$ , the number  $N_{t_0} \times N_\tau$  of grid points and the number of individual summation steps  $N_{\text{sums}}$ . However, for the timings presented in figure 1, we are mostly in regimes where the dominant scaling for `rect` windows is with  $N_{t_0} \times N_\tau$  and the cost for `exp` windows is dominated by the  $N_{\text{sums}}$  scaling. Hence, the results are still quite consistent with the simpler model originally introduced in [1].

After fitting the more general model to the measured timings, the fit coefficients can be approximately converted to the PGM per-machine ‘timing constants’  $\mathbf{c}_r$  and  $\mathbf{c}_e$  for the two window functions. In the original timing model, these are interpreted directly as the cost to compute the (weighted) sums over atoms at each step, but under the more general model, each corresponds to a combination of the three contributions. Results are listed in table 1.

For fairness in these comparisons, we note that the current `LALPulsar` CPU code may not be fully optimised. Two algorithmic optimisations are already included in it: (i) the summation truncation for `rect` windows discussed above, and (ii) a lookup-table (LUT) based ‘fast exponential’ function to at least somewhat reduce the cost of the



**Table 1.** Timing results from figure 1 converted to the timing constants  $\tau_r$ ,  $\tau_e$  as introduced by [1], approximately corresponding to the cost of each individual summation step. First the timings are fit with the extended model from Appendix A (fit errors are  $< 1\%$ ) and, assuming dominant  $N_{t_0} \times N_\tau$  scaling for rectangular windows and  $N_{\text{sums}}$  scaling for exponential windows, the coefficients are converted to  $\tau_r$  and  $\tau_e$ . Appendix A has additional details and complementary example fits of the more general timing model to different scaling regimes. Note that  $\tau_e < \tau_r$  for the GPUs does not mean that the overall search for an exponential window is faster than for a rectangular window, as these constants are multiplied with different summation counters, see (A.1) and (A.2).

CPU/GPU	$\tau_r$ [s]	$\tau_e$ [s]	
Core2Duo 2.6 GHz	$4.2 \cdot 10^{-8}$	$1.3 \cdot 10^{-7}$	from [1]
i5-6200U	$6.4 \cdot 10^{-8}$	$1.1 \cdot 10^{-7}$	
Xeon X5675	$3.5 \cdot 10^{-8}$	$7.0 \cdot 10^{-8}$	
Xeon E5-2630	$3.4 \cdot 10^{-8}$	$5.7 \cdot 10^{-8}$	
Xeon E5-2650	$3.6 \cdot 10^{-8}$	$6.0 \cdot 10^{-8}$	
GTX-1050	$6.1 \cdot 10^{-9}$	$1.4 \cdot 10^{-9}$	
GTX-1060	$4.8 \cdot 10^{-9}$	$9.1 \cdot 10^{-10}$	
GTX-1070	$4.2 \cdot 10^{-9}$	$7.3 \cdot 10^{-10}$	
GTX-1080	$4.4 \cdot 10^{-9}$	$6.2 \cdot 10^{-10}$	
Tesla-V100	$4.3 \cdot 10^{-9}$	$4.6 \cdot 10^{-11}$	

exp window (see also subsection 4.3 below). While the current code does not make full use of the theoretical peak floating-point operations per second (flops) of modern CPUs, the atoms-based transient- $\mathcal{F}$ -statistic algorithm fundamentally has a significant computing cost contribution from index comparisons and other non-floating-point operations, so that it is not currently obvious if any significant further optimisations would be feasible. Meanwhile, the main finding of these tests is that already a *straightforward* pyCUDA port of the algorithm, without additional optimisations, yields significant speed-ups.

#### 4.2. Memory

GPU applications are often memory-limited. However, for the transient  $\mathcal{F}$ -statistic map, we do not expect GPU memory to be a significant constraint, as we see in the following. With the current approach, the input atoms need to be transferred to GPU memory only for a single  $\lambda$  parameter space point at a time, then the  $\mathcal{F}_{mn}(\lambda)$  matrix is computed and returned. Hence, the peak GPU memory usage of input plus output matrices is expected to be

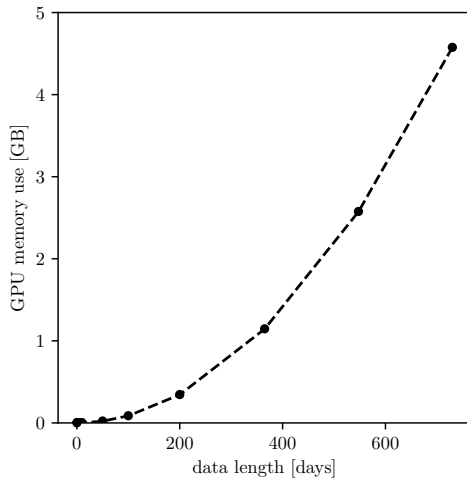
$$M[\text{bytes}] = 4 ( 7N_{\text{SFT}} + N_{t_0} N_\tau ), \quad (9)$$

where 4 bytes is the base size of a `real32` number in the underlying NumPy [23] package. While the input array size grows only linearly with  $N_{\text{SFT}}$ , assuming  $dt_0 = d\tau = T_{\text{SFT}}$  the  $\mathcal{F}_{mn}$  matrix grows quadratically and will dominate memory usage at long  $T_{\text{obs}}$ . However, in practice one might want to choose an undersampling of  $t_0, \tau$ .

A comparison of this expectation with practical memory usage measurements is presented in figure 2. For  $T_{\text{SFT}} = 1800\text{s}$  and  $dt_0 = d\tau = T_{\text{SFT}}$ , the memory usage reaches only about 1.1 GB for a year of data, and with undersampling even much longer data sets would remain easily feasible on current GPUs, even when multiple jobs need to run on a single device.

GPU transient  $\mathcal{F}$ -statistic

8

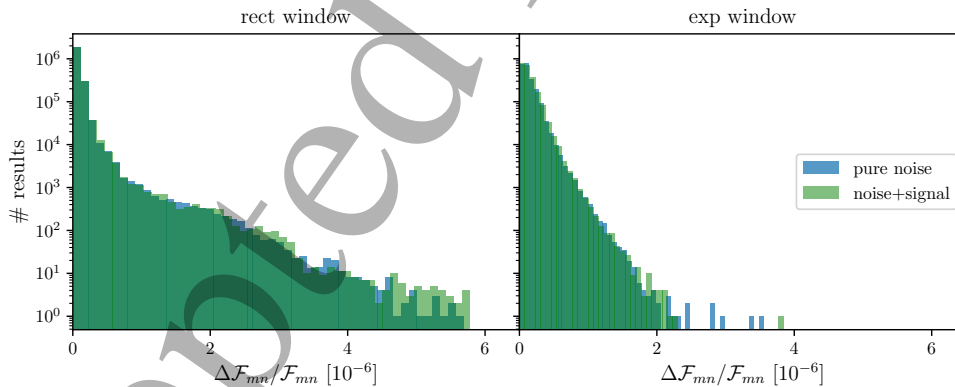


**Figure 2.** GPU memory usage on a GeForce GTX 1070 (8 GB RAM) with CUDA V8.0.61. Measured with  $T_{\text{SFT}} = 1800$  s and a resolution of  $dt_0 = d\tau = T_{\text{SFT}}$  at all  $T_{\text{obs}}$ . Each data point is the difference between the output of a call to `pycuda.driver.mem_get_info()` right before allocating input and output arrays with `pycuda.gparray`, and a call right afterwards. The dashed line is the expected  $4(7N_{\text{SFT}} + N_{t_0}N_\tau)$  scaling (in bytes). As  $T_{\text{obs}} \rightarrow 0$ , we find that the base memory use for the kernel itself (and any other possible overheads) seems to be only about 2–4 MB.

## 4.3. Accuracy

The original LALPulsar implementation is already using single precision for the atoms and the  $\mathcal{F}$ -statistic itself, so in contrast to some other GPU use cases [24] it was not necessary to reduce the code’s internal precision for the pyCUDA version. However, the  $\mathcal{F}$ -statistic algorithm is already known to produce slightly different numerical results on different CPU platforms, so it is worth checking the typical amount of differences in the transient  $\mathcal{F}$ -statistic between LALPulsar and pyCUDA versions.

As demonstrated for a particular test case in figure 3, we typically find negligibly small differences, not larger than other implementation- and platform-dependent



**Figure 3.** Comparison of  $\mathcal{F}_{mn}$  computed with the LALPulsar and pyCUDA implementations. Each histogram gives the differences  $\Delta\mathcal{F}_{mn}$  between the two implementations for a certain transient window, and either for pure Gaussian noise or also including a (confidently detectable,  $\max 2\mathcal{F} \approx 263$ ) signal injection with matching window function. The histograms are taken over all individual  $\mathcal{F}_{mn}$  values for 1000 frequency bins over a 1-day data set with  $T_{\text{SFT}} = 1800$ . The GPU for this test was a GeForce GTX 1070.

**Table 2.** Example per-template timings of the  $\mathcal{F}_{mn}$ -map cost for post-glitch search setups with  $\tau \in [2T_{\text{SFT}}, T_{\text{obs}} = 4 \text{ months}]$  and  $t_0$  in the band listed. Also listed are GPU-over-CPU speedup factors and extrapolated total costs for example searches with template counts  $N_\lambda$  from [26, 27]. The CPU columns are for a single core of an Intel Xeon E5-2650 and the GPU columns for an Nvidia Tesla-V100. The per-template timings are averaged over 5 runs of 100–1000 templates each, except for the last CPU case (only 2 single-template runs).

$t_0$ band (days)	window	seconds per template		GPU speedup	Vela: $N_\lambda \approx 2.2 \cdot 10^6$		Crab: $N_\lambda \approx 1.7 \cdot 10^8$		
		CPU	GPU		CPU	GPU	CPU	GPU	
0	rect	$3.5 \cdot 10^{-4}$	$2.4 \cdot 10^{-2}$	< 1	—	not dominated by $\mathcal{F}_{mn}$ cost			
0	exp	$5.7 \cdot 10^{-1}$	$2.3 \cdot 10^{-2}$	$\approx 25$	—	not dominated by $\mathcal{F}_{mn}$ cost			
1	rect	$1.1 \cdot 10^{-2}$	$2.6 \cdot 10^{-2}$	< 1	0.3 d	0.7 d	21 d	51 d	
1	exp	$2.8 \cdot 10^1$	$3.7 \cdot 10^{-2}$	$\approx 750$	2 y	1 d	150 y	72 d	
120	rect	$1.4 \cdot 10^0$	$1.6 \cdot 10^{-1}$	$\approx 9$	35 d	4 d	7 y	1 y	
120	exp	$1.7 \cdot 10^3$	$7.7 \cdot 10^{-1}$	$\approx 2300$	120 y	20 d	9300 y	4 y	

variations in the  $\mathcal{F}$ -statistic known from other work (e.g. [25]).

One relevant implementation detail of the `LALPulsar` code is the previously mentioned lookup-table (LUT) based ‘fast exponential’ function. This can actually lead to differences with `pyCUDA` for exponential windows of up to  $\sim 10\%$ , but figure 3 shows results after replacing it with the `exp()` function of the C standard library, thus verifying that the difference did not come from a loss of accuracy with the new `pyCUDA` implementation.

## 5. Conclusion and applications

The significant speedup achieved with our `pyCUDA` implementation of the transient  $\mathcal{F}$ -statistic will allow for a wider scope of searches for long-duration transient GWs. We now discuss a few example applications that would be hard resource-wise, or even prohibitive, on CPUs but could become viable with GPUs.

Let us first consider the natural use case of a GW data analysis triggered by electromagnetic (EM) observations of a pulsar glitch. Quasi-monochromatic GW emission, which the  $\mathcal{F}$ -statistic is sensitive to, could be associated with the post-glitch relaxation phase. Depending on the pulsar, this can have timescales of days to months [11, 12] with an exponential decay, guiding the parameter range for a GW transient search and motivating an exponential amplitude window.

As simple examples for possible transient search setups, assume we look at  $T_{\text{obs}} = 4 \text{ months}$  of data, signal durations  $\tau \in [2T_{\text{SFT}}, T_{\text{obs}}]$ , and for the starting time  $t_0$  we choose either (i)  $t_0 \in [0, 1 \text{ day}]$  or (ii)  $t_0 \in [0, T_{\text{obs}} - 2T_{\text{SFT}}]$ , with grid steps equal to  $T_{\text{SFT}} = 1800 \text{ s}$  for both  $t_0$  and  $\tau$ . Setup (i) represents some uncertainty in GW signal start time in relation to the glitch, while (ii) would instead be a completely generic transient search for any signals after the glitch, not necessarily correlated with the glitch time. Another possible setup, with a single fixed  $t_0$ , would be appropriate if models predict unequivocally that the GW emission starts at the same time as the glitch; it is computationally cheap enough to not consider in detail here.

With these parameters, we obtain the per-template runtimes listed in Table 2. Comparing a Tesla-V100 GPU with single-core runs on a Xeon E5-2650 CPU, we find that for rectangular transient windows the GPU is actually still slower than the CPU

in case (i) and only outperforms it by a factor of  $\sim 10$  for case (ii). But for exponential windows, the Tesla beats the Xeon significantly by factors of approximately 750 and 2300. Still, with a single GW search template matching the post-glitch radio timing solution (at  $f_{\text{GW}} = 2f_{\text{spin}}$ ), any of these analyses would be computationally trivial even on a single CPU.

However, it is reasonable to allow for some mismatch between the radio timing and GW frequency evolution due to the perturbed state of the NS after a glitch. For comparison, the ‘narrow-band’ search for CWs from known pulsars in the first aLIGO run [26] (using 121 days of data, matching the 4 months considered in the Table 2 timings) covered some ranges in frequency  $f$  and spindown  $\dot{f}$  for each of its 11 targets, with totals of e.g.  $N_\lambda \approx 2.2 \cdot 10^6$  templates for the Vela pulsar and  $N_\lambda \approx 1.7 \cdot 10^8$  for the Crab pulsar.¶ (These values have been corrected in Table I of an erratum [27] to [26].)

Multiplying these numbers of templates with the per-template transient  $\mathcal{F}$ -statistic costs (which in these setups again dominates over the rest of the search algorithm), we find that in the more expensive case (ii) ( $t_0 \in [0, T_{\text{obs}} - 2T_{\text{SFT}}]$ ), a single Tesla could perform an exponential-window transient analysis over the Vela band in less than 3 weeks, while the same analysis would take 120 years on a single Xeon-E5-like CPU core; or equivalently would require about 2300 CPU cores to only take the same 3 weeks as the single Tesla. For the wider Crab range (which was chosen in [26] to account for the Crab’s strong spindown), even the Tesla would still need 4 years. The more relevant scenario in practice is that of a moderate  $t_0$  range, e.g. case (i) with a 1-day range. Such a search over the Vela band would take a single day on a Tesla, while requiring over 700 CPU core-days. Even the wider Crab range (chosen in [26] to account for the Crab’s stronger spindown) would become feasible with 2–3 months on a single Tesla, instead of occupying a significant part of a large CPU cluster for several weeks.

In summary, performing routine transient  $\mathcal{F}$ -statistic analyses of *all* observed glitches in known galactic pulsars during a GW observation run – with reasonably wide bands in  $f$  and  $\dot{f}$  (similar to those used in [26], or only slightly reduced) and in  $t_0$  – becomes feasible with a few dedicated GPU systems. Choosing GPUs over CPUs for this kind of search is also efficient in terms of actual money cost, when we compare the cost of adding a single GPU to an existing computer or adding another CPU-only worker node to a cluster. For example, the Tesla-V100 on the CIT cluster used for the above benchmarks cost about 6200 USD, while a single Xeon E5-2650 processor (8 cores) was initially listed by Intel at 1100 USD. Allowing for a modest overhead for other components in a CPU cluster node, a Tesla would thus already be efficient in purchasing cost if faster than about  $4 \cdot 8 = 32$  CPU cores. With power consumptions of 250 W for the Tesla and  $\sim 95$  W (Intel’s listed Thermal Design Power) for a Xeon under full load, the electricity cost comparison is even more favourable. Numbers will of course vary for different hardware and purchase times, with more modern devices typically delivering more power per USD and per W.

Similar estimates as for EM-triggered glitch searches also apply when considering the follow-up [14, 30, 31] of significant or marginal detection candidates produced by wide-parameter space CW searches [3]. Even though those searches target perfectly persistent signals, they can also produce candidates if there are sufficiently strong

¶ Both are interesting targets for post-glitch transient searches: Vela last glitched on 2016-12-12 [28], during the Advanced LIGO O2 run, and the Crab last glitched on 2017-11-08 [29], unfortunately after the end of O2.

transient events in the data [32]. A comprehensive transient-aware follow-up, with the goal of either verifying the presumed persistent nature or uncovering a transient signal instead, needs to not only target the exact phase-evolution parameters  $\lambda$  of the candidate, but search a wider band around it to account for degeneracies with the transient evolution parameters. Reducing the computational cost of each candidate's follow-up directly translates into a larger number of candidates that can be analysed, so that the overall threshold of the CW search can be lowered and a better search sensitivity can be achieved.

The data length and  $\{t_0, \tau\}$  ranges in this scenario can be longer than in the EM-triggered post-glitch scenario: the aLIGO runs O1–O3 took data for / are scheduled for 4, 9 and 12 months respectively [33], and for the follow-up of a strong candidate data from multiple observing runs could get combined. The range of phase evolution parameters  $\lambda$  that should be searched for full coverage depends on the exact setup of the CW search and on possible intermediate follow-up steps; but the scaling of the transient  $\mathcal{F}$ -statistic cost would be similar to the more expensive full- $t_0$ -range case (ii) considered in the EM-triggered example above (see Appendix A for the full timing model) and the improvements in accessible search volume using a small number of GPUs over CPUs will be at least similar.

In the longer term, untriggered all-sky searches for long-duration transients are of high interest. Similarly to all-sky CW searches, they have the potential to discover a population of electromagnetically dark NSs, for example glitching pulsars with their beam pointed away from Earth. The sensitivity of all-sky searches is directly limited by how densely they can cover the  $\lambda$  parameter space at a fixed computational budget. [34, 35]. Hence, adding transient parameters at first significantly reduces the overall sensitivity of a blind search. But speeding up the transient part by orders of magnitude could still make a combined search for CWs and transients feasible in the long run, when large numbers of specialised GPUs become available in high-performance clusters. This could also be of interest for distributed volunteer computing [Einstein@Home 36]: the speedup with consumer GPUs (e.g. the GTX family) is more modest, but still significant. In practice, though, the more promising approach for blind transient searches might be to apply a cheap add-on transient modification, like that introduced in [32], to a semi-coherent CW algorithm as a first search stage, then apply the fully-coherent transient  $\mathcal{F}$ -statistic only in a follow-up step.

In any of these scenarios, while we have focussed on the fact that the pyCUDA version can bring down the cost of exponential-windowed transients significantly, the cost for rectangular windows always remains smaller, so that in practice whenever exponential windows are feasible, it is also cheap and natural to run *both* analyses and evaluate a posteriori which one fits the data better. Different window functions for the amplitude evolution could also be considered, and would generically follow the GPU kernel grid setup and timing model for the exponential window, since it does not assume any function-specific optimisations.

### Acknowledgments

The authors would like to thank Reinhard Prix for feedback on the manuscript as well as him and Karl Wette for development and continued support of the underlying LALSuite  $\mathcal{F}$ -statistic code. Thanks to Chris Messenger and Stuart Anderson for support with the GPU test systems at Glasgow and CIT. We also thank the CQG

referees for thoughtful suggestions on improving the paper, especially on extending the discussion of efficiency, practical applications and costs.

DK was funded under the EU Horizon2020 framework through the Marie Skłodowska-Curie grant agreement 704094 GRANITE, and would also like to acknowledge the Horizon2020 ASTERICS-OBELICS International School (Annecy 2017) that inspired this investigation into pyCUDA and GPU applications.

## Appendix

### Appendix A. Generalising the PGM2011 timing model

Here we revisit the timing model for computing  $\mathcal{F}_{mn}$  maps introduced in Appendix A3 of [1]. Their equations (A13) and (A14) give the computing cost for a single- $\lambda$   $\mathcal{F}_{mn}$  map with either exponential

$$\mathbf{c}_{\mathcal{F}\text{map}}^e \approx \mathbf{c}_e \frac{\Delta t_0}{dt_0} \frac{\Delta \tau}{d\tau} \frac{(\tau_{\min} + \Delta\tau/2)}{T_{\text{SFT}}} \approx \mathbf{c}_e N_{t_0} N_\tau \frac{(\tau_{\min} + \Delta\tau/2)}{T_{\text{SFT}}} \quad (\text{A.1})$$

or rectangular window functions:

$$\mathbf{c}_{\mathcal{F}\text{map}}^r = \mathbf{c}_r \frac{\Delta t_0}{dt_0} \frac{(\tau_{\min} + \Delta\tau)}{T_{\text{SFT}}} \approx \mathbf{c}_r N_{t_0} \frac{(\tau_{\min} + \Delta\tau)}{T_{\text{SFT}}} = \mathbf{c}_r N_{t_0} \frac{\tau_{\max}}{T_{\text{SFT}}}. \quad (\text{A.2})$$

The timing constants  $\mathbf{c}_e$  and  $\mathbf{c}_r$  are interpreted as the cost to compute the (weighted) sums over atoms at each step. The exponential model corresponds to a ‘generic’ case where all quantities have to be re-evaluated at each step, while the rectangular case reuses partial sums as discussed before.

We note now that this formulation of the timing model does not explicitly include the cost of computing the antenna pattern matrix determinant  $\hat{D}$  and the  $\mathcal{F}$ -statistic itself, which is done once for each  $(m, n)$  pair after all sums have been computed and hence is independent of the window function choice.<sup>+</sup> We can include this contribution by adding a term  $+\mathbf{c}_{\mathcal{F}} N_{t_0} N_\tau$  to both cases. It will be very subdominant for exponential windows, where the summations term grows much faster than  $N_{t_0} N_\tau$ , but can be relevant for rectangular windows where the summations term is more efficient.

Another small contribution to the timing model is from setup and index-lookup costs that scale with the total number of SFTs handed to the  $\mathcal{F}$ -statistic-map function; for completeness we include a common term  $\mathbf{c}_{\text{SFTs}} N_{\text{SFT}}$ .

In addition, Eqs. (A.1) and (A.2) only hold true if the full range of transient signal durations explored by the  $\mathcal{F}_{mn}$  map is fully contained within the available data range, that is when  $t_{0\text{max}} + \tau_{\text{max}} < T_0 + T_{\text{data}}$ . (We call this the ‘embedded’ case below.) Otherwise, i.e. if some of the transient windows overlap the end of the available data, by convention the LALPulsar code still returns results for the full rectangular  $\mathcal{F}_{mn}$  matrix, but truncates the atoms summations. Thus, the total computing cost in such cases is lower than estimated by Eqs. (A.1), (A.2) and using them to fit the timing constants from runtime measurements as in Sec. 4.1 would yield inconsistent results.

Hence, we generalise the PGM timing model by introducing  $N_{\text{sums}}$  as the effective number of summation steps for an  $\mathcal{F}_{mn}$  map, which depends on the window type,  $T_{\text{data}}$ , and the ranges of both  $t_0$  and  $\tau$ :

$$\mathbf{c}_{\mathcal{F}\text{map}}^e \approx \mathbf{c}_{\text{SFTs}} N_{\text{SFT}} + \mathbf{c}_{\text{sums}}^e N_{\text{sums}}^e + \mathbf{c}_{\mathcal{F}} N_{t_0} N_\tau, \quad (\text{A.3})$$

<sup>+</sup> In its scalings, this extra cost is degenerate with the marginalisation cost  $\mathbf{c}_{\text{marg}}$  of PGM’s Eq. (15), in search code executions where both  $\mathcal{F}_{mn}$  maps and marginal Bayes factors are computed; so it was effectively included in PGM’s overall code timing, but just attributed to a different part of the model.

GPU transient  $\mathcal{F}$ -statistic

13

$$\mathbf{c}_{\mathcal{F}\text{map}}^r \approx \mathbf{c}_{\text{SFTs}} N_{\text{SFT}} + \mathbf{c}_{\text{sums}}^r N_{\text{sums}}^r + \mathbf{c}_{\mathcal{F}} N_{t_0} N_{\tau}. \quad (\text{A.4})$$

For rectangular windows, we have

$$N_{\text{sums}}^r = \sum_{m=1}^{N_{t_0}} \frac{\min(T_{\text{data}} - t_{0m}, \tau_{\text{max}})}{T_{\text{SFT}}}, \quad (\text{A.5})$$

which reduces to PGM's  $N_{\text{sums}}^r = N_{t_0} \tau_{\text{max}} / T_{\text{SFT}}$  in the special 'embedded' case that PGM considered, and to  $N_{\text{sums}}^r = 0.5 N_{t_0} \tau_{\text{max}} / T_{\text{SFT}}$  in the special case of  $N_{t_0} dt_0 = N_{\tau} d\tau = T_{\text{data}} - 2T_{\text{SFT}}$  that we used for the timing results in Sec. 4.1.

For exponential windows, we also need to note that the current code's convention, as introduced in Eq. (18) of [1], is that an exponential window with duration parameter  $\tau$  covers an effective length of  $3\tau / T_{\text{SFT}}$  atoms. (The exponential decay is not cut off after only one, but after three e-folds, where the remaining SNR would be much more negligible.) Hence, PGM's original timing constant  $\mathbf{c}_e$  effectively contains a factor of 3 (from counting all steps in  $\tau$ ) that we now include in  $N_{\text{sums}}^e$  instead:

$$N_{\text{sums}}^e = \sum_{m=1}^{N_{t_0}} \sum_{n=1}^{N_{\tau}} \frac{\min(T_{\text{data}} - t_{0m}, 3\tau_n)}{T_{\text{SFT}}}. \quad (\text{A.6})$$

In the 'embedded' case this reduces to  $3N_{t_0} \sum_{n=1}^{N_{\tau}} \tau_n / T_{\text{SFT}} = 3N_{t_0} N_{\tau} (\tau_{\text{min}} + 0.5\Delta\tau) / T_{\text{SFT}}$ , equivalent to PGM's result up to the factor of 3.

In practice, on each architecture we can use these more general equations (A.3)–(A.6) to fit the four timing constants  $\{\mathbf{c}_{\text{SFTs}}, \mathbf{c}_{\mathcal{F}}, \mathbf{c}_{\text{sums}}^r, \mathbf{c}_{\text{sums}}^e\}$  from a variety of setups (in terms of  $T_{\text{data}}, [t_{0\text{min}}, t_{0\text{max}}], [\tau_{\text{min}}, \tau_{\text{max}}]$ ), then consider the special 'embedded' case\* (and  $N_{\tau} \gg 1, \tau_{\text{max}} \gg \tau_{\text{min}}$ ) to directly compare to [1] by

$$\mathbf{c}_r = \mathbf{c}_{\text{sums}}^r + \frac{T_{\text{SFT}}}{\tau_{\text{max}}} \left( \mathbf{c}_{\text{SFTs}} \frac{N_{\text{SFT}}}{N_{t_0}} + \mathbf{c}_{\mathcal{F}} N_{\tau} \right) \approx \mathbf{c}_{\text{sums}}^r + \mathbf{c}_{\mathcal{F}}, \quad (\text{A.7})$$

$$\mathbf{c}_e = 3\mathbf{c}_{\text{sums}}^e + \frac{T_{\text{SFT}}}{\tau_{\text{min}} + 0.5\Delta\tau} \left( \mathbf{c}_{\text{SFTs}} \frac{N_{\text{SFT}}}{N_{t_0} N_{\tau}} + \mathbf{c}_{\mathcal{F}} \right) \approx 3\mathbf{c}_{\text{sums}}^e + \frac{2}{N_{\tau}} \mathbf{c}_{\mathcal{F}} \approx 3\mathbf{c}_{\text{sums}}^e. \quad (\text{A.8})$$

Using a set of timing runs that in addition to those in section 4.1 also cover many different combinations of  $\{T_{\text{data}}, N_{t_0}, N_{\tau}\}$ , and also measuring *only* the execution time of the actual  $\mathcal{F}$ -statistic map function (while in section 4.1 the whole search call is timed, including the contribution of computing the atoms, which is usually subdominant but not in the limit of low  $N_{t_0} N_{\tau}$  and  $N_{\text{sums}}$ ), we do a detailed fit of the full timing model of (A.3) and (A.4), in the following iterative steps to ensure convergence:

- (i) fit the  $\mathbf{c}_{\text{SFTs}} N_{\text{SFT}}$  term only to short data sets with  $N_{\text{sums}} \leq 100$
- (ii) using this fixed  $\mathbf{c}_{\text{SFTs}}$ , fit  $\mathbf{c}_{\text{sums}}^r N_{\text{sums}}^r + \mathbf{c}_{\mathcal{F}} N_{t_0} N_{\tau}$  for rectangular windows
- (iii) using fixed  $\mathbf{c}_{\text{SFTs}}$  and  $\mathbf{c}_{\mathcal{F}}$ , fit  $\mathbf{c}_{\text{sums}}^e N_{\text{sums}}^e$  for exponential windows

\* The example used for timing in Appendix A3 of [1] is 1 year of data with  $\tau \in [0.5, 14.5]$  days; which means that with  $t_0$  close to the end of the year and  $\tau_{\text{max}} = 14.5$  days the overlap is at most 4% and the deviations from the fully-embedded special case used for this comparison are smaller than typical timing uncertainties.

## REFERENCES

14

We find e.g.

$$\mathbf{c}_{\mathcal{F}_{\text{map}}}^r \approx ((2.80 \pm 0.03)N_{\text{SFT}} + (0.96 \pm 0.08)N_{\text{sums}}^r + (5.59 \pm 0.06)N_{t_0}N_{\tau}) 10^{-8} \text{ s} \quad (\text{A.9})$$

$$\mathbf{c}_{\mathcal{F}_{\text{map}}}^e \approx ((2.80 \pm 0.03)N_{\text{SFT}} + (3.55 \pm 0.03)N_{\text{sums}}^e + (5.59 \pm 0.06)N_{t_0}N_{\tau}) 10^{-8} \text{ s} \quad (\text{A.10})$$

for the i5-6200U laptop CPU (corresponding to PGM constants  $\mathbf{c}_r = (6.36 \pm 0.08)10^{-8}$  and  $\mathbf{c}_e = (1.07 \pm 0.01)10^{-7}$ ); and

$$\mathbf{c}_{\mathcal{F}_{\text{map}}}^r \approx ((2.59 \pm 0.02)N_{\text{SFT}} + (0.27 \pm 0.02)N_{\text{sums}}^r + (3.09 \pm 0.02)N_{t_0}N_{\tau}) 10^{-8} \text{ s} \quad (\text{A.11})$$

$$\mathbf{c}_{\mathcal{F}_{\text{map}}}^e \approx ((2.59 \pm 0.02)N_{\text{SFT}} + (2.22 \pm 0.03)N_{\text{sums}}^e + (3.09 \pm 0.02)N_{t_0}N_{\tau}) 10^{-8} \text{ s} \quad (\text{A.12})$$

for the Xeon X5675 workstation CPU (corresponding to PGM constants  $\mathbf{c}_r = (3.26 \pm 0.02)10^{-8}$  and  $\mathbf{c}_e = (6.67 \pm 0.08)10^{-8}$ ). These results agree reasonably well with those obtained on the same systems, but with fixed  $N_{t_0}, N_{\tau}$  in relation to  $T_{\text{data}}$  and with simplified fits, as presented in table 1. While the error bars from fitting alone appear too small to explain the remaining differences of 0.5–7%, it is likely that variations in system configuration and load between timing runs are the main culprit.

## References

- [1] Prix R, Giampanis S and Messenger C 2011 *Phys. Rev. D* **84** 023007 [arXiv:1104.1704]
- [2] Prix R (for the LIGO Scientific Collaboration) 2009 *Gravitational Waves from Spinning Neutron Stars (Astrophys. Space Sci. Lib. vol 357)* (Springer Berlin Heidelberg) chap 24, pp 651–685 ISBN 978-3-540-76964-4 URL <https://dcc.ligo.org/LIGO-P060039/public>
- [3] Riles K 2017 *Mod. Phys. Lett. A* **32** 1730035 [arXiv:1712.05897]
- [4] Aasi J *et al.* (LIGO Scientific Collaboration) 2015 *Class. Quant. Grav.* **32** 074001 [arXiv:1411.4547]
- [5] Acernese F *et al.* (Virgo Collaboration) 2015 *Class. Quant. Grav.* **32** 024001 [arXiv:1408.3978]
- [6] Jaranowski P, Królak A and Schutz B F 1998 *Phys. Rev. D* **58** 063001 [arXiv:gr-qc/9804014]
- [7] Cutler C and Schutz B F 2005 *Phys. Rev. D* **72** 063006 [arXiv:gr-qc/0504011]
- [8] Aasi J *et al.* (LIGO Scientific Collaboration and Virgo Collaboration) 2015 *Astrophys. J.* **813** 39 [arXiv:1412.5942]
- [9] Zhu S J *et al.* 2016 *Phys. Rev. D* **94** 082008 [arXiv:1608.07589]
- [10] Abbott B P *et al.* (LIGO Scientific Collaboration and Virgo Collaboration) 2017 *Phys. Rev. D* **96** 122004 [arXiv:1707.02669]
- [11] Lyne A G, Shemar S L and Smith F G 2000 *Mon. Not. R. Astron. Soc.* **315** 534–542
- [12] Haskell B and Antonopoulou D 2014 *Mon. Not. Roy. Astron. Soc.* **438** 16 [arXiv:1306.5214]
- [13] Klöckner A, Pinto N, Lee Y, Catanzaro B, Ivanov P and Fasih A 2012 *Parallel Computing* **38** 157–174 ISSN 0167-8191
- [14] Ashton G and Prix R 2018 *Phys. Rev. D* **97** 103020 [arXiv:1802.05450]
- [15] Ashton G and Keitel D 2018 Pyfstat-v1.2 URL <https://doi.org/10.5281/zenodo.1243931>
- [16] LSC Algorithm Library - LALSuite (free software) URL <https://git.ligo.org/lscsoft/lalsuite>



## REFERENCES

15

- [17] Williams P R and Schutz B F 1999 *AIP Conf. Proc.* **523** 473 [arXiv:gr-qc/9912029]
- [18] Prix R 2009 *The F-statistic and its implementation in ComputeFStatistic\_v2* Tech. Rep. LIGO-T0900149-v6 last updated 2018 URL <https://dcc.ligo.org/LIGO-T0900149/public>
- [19] Prix R and Krishnan B 2009 *Class. Quant. Grav.* **26** 204013 [arXiv:0907.2569]
- [20] Keitel D, Prix R, Papa M A, Leaci P and Siddiqi M 2014 *Phys. Rev. D* **89** 064023 [arXiv:1311.5738]
- [21] Beazley D M 1996 SWIG: An Easy to Use Tool for Integrating Scripting Languages with C and C++ *Proc. 4th Conf. USENIX Tcl/Tk Workshop* (Berkeley, CA, USA: USENIX Association) pp 15–15 URL <http://dl.acm.org/citation.cfm?id=1267498.1267513>
- [22] Beazley D M *et al.* SWIG - Simplified Wrapper and Interface Generator URL [www.swig.org](http://www.swig.org)
- [23] Oliphant T E 2006 *A guide to NumPy* (Trelgol Publishing)
- [24] Navarro C A, Hitschfeld-Kahler N and Mateu L 2014 *Communications in Computational Physics* **15** 285329
- [25] Prix R 2011 *F-statistic bias due to noise-estimator* Tech. Rep. LIGO-T1100551 URL <https://dcc.ligo.org/LIGO-T1100551/public>
- [26] Abbott B P *et al.* (Virgo, LIGO Scientific) 2017 *Phys. Rev. D* **96** 122006 [arXiv:1710.02327]
- [27] Abbott B P *et al.* (Virgo, LIGO Scientific) 2018 *Phys. Rev. D* **97** 129903
- [28] Palfreyman J, Dickey J M, Hotan A, Ellingsen S and van Straten W 2018 *Nature* **556** 219–222
- [29] Shaw B *et al.* 2018 *Mon. Not. R. Astron. Soc.* **478** 3832–3840 [arXiv:1805.05110]
- [30] Shaltev M and Prix R 2013 *Phys. Rev. D* **87** 084057 [arXiv:1303.2471]
- [31] Papa M A *et al.* 2016 *Phys. Rev. D* **94** 122006 [arXiv:1608.08928]
- [32] Keitel D 2016 *Phys. Rev. D* **93** 084024 [arXiv:1509.02398]
- [33] Abbott B P *et al.* (VIRGO, LIGO Scientific) 2018 *Living Rev. Relativity* **21:3** [arXiv:1304.0670] URL <https://link.springer.com/article/10.1007/s41114-018-0012-9>
- [34] Prix R and Shaltev M 2012 *Phys. Rev. D* **85** 084010 [arXiv:1201.4321]
- [35] Wette K 2012 *Phys. Rev. D* **85** 042003 [arXiv:1111.5650]
- [36] Allen B *et al.* Einstein@Home distributed computing project URL <https://einsteinathome.org>