

# Improving the prospects for multimessenger astronomy with early-warning detection of compact binary coalescence.

Kipp Cannon<sup>1</sup>, Romain Cariou<sup>2</sup>, Adrian Chapman<sup>3</sup>, M. Crispin-Ortuzar<sup>4</sup>, Nickolas Fotopoulos<sup>3</sup>, Melissa Frei<sup>5</sup>, Chad Hanna<sup>6</sup>, Erin Kara<sup>7</sup>, Drew Keppel<sup>8,9</sup>, Laura Liao<sup>10</sup>, Stephen Privitera<sup>3</sup>, Antony C. Searle<sup>3</sup>, Leo Singer<sup>3</sup>, Alan J. Weinstein<sup>3</sup>

<sup>1</sup> Canadian Institute for Theoretical Astrophysics, Toronto, ON, Canada

<sup>2</sup> Département de physique, École Normale Supérieure de Cachan, 61 Avenue du Président Wilson, 94235 Cachan Cedex, France

<sup>3</sup> LIGO Laboratory - California Institute of Technology, Pasadena, CA, USA

<sup>4</sup> Facultat de Física, Universitat de València, E-46100 Burjassot, Spain

<sup>5</sup> The University of Texas at Austin, Austin, TX, USA

<sup>6</sup> Perimeter Institute for Theoretical Physics, Waterloo, ON, Canada

<sup>7</sup> Department of Physics and Astronomy, Barnard College, Columbia University, New York, NY 10027, USA

<sup>8</sup> Albert-Einstein-Institut, Max-Planck-Institut für Gravitationsphysik, Hannover, Germany

<sup>9</sup> Leibniz Universität Hannover, Hannover, Germany

<sup>10</sup> Ryerson University, Toronto, ON, Canada

**Abstract.** The rapid detection of compact binary coalescence with a network of advanced gravitational-wave detectors will offer a unique opportunity for multimessenger astronomy. Prompt detection alerts to the astronomical community may make it possible to observe the onset of electromagnetic emission from compact binary coalescence. We demonstrate a computationally practical analysis strategy that produces early warning triggers even before gravitational radiation from the final merger has arrived at the detectors. With current rate estimates for the Advanced LIGO design configuration, we should detect  $\sim 10$  sources earlier than 10 seconds before merger in 1 year of livetime.

PACS numbers: 95.55.Ym, 84.30.Vn, 95.75.Wx

## 1. Introduction

The coalescence of compact binary systems consisting of neutron stars (NS) and/or black holes (BH) is the most promising source of gravitational radiation for Advanced LIGO [1], Virgo [2], GEO [3], and LCGT [4]. Tens of binary coalescence events are expected to be observed in the advanced detector era later this decade [5].

As a compact binary system loses energy to gravitational waves, its orbital separation decreases. This causes a run-away inspiral with the gravitational-wave amplitude and frequency increasing until the system eventually merges near the innermost stable circular orbit (ISCO). If a neutron star is involved it may

*Would the introduction be more effective without this paragraph?*

become tidally disrupted near the merger. This disrupted matter can fuel a bright electromagnetic counterpart in the system’s final moments as a binary [6].

Prompt electromagnetic emission can arise as shells of relativistically outflowing matter collide in the inner shock. Such an inner shock from a compact binary coalescence is believed to be a mechanism for short gamma-ray bursts (short GRBs) [7, 8]. The same inner shocks, or potentially reverse shocks, can produce a bright accompanying optical flash [9]. Prompt emission is a probe into the extreme initial conditions of the outflow, in contrast with afterglows, which are relatively insensitive to initial conditions. Optical flashes have only been observed for a handful of long GRBs [10, 11, 12, 13, 14, 15] by telescopes with extremely rapid response or, in the case of GRB 080319B, by pure serendipity, where several telescopes were observing a previous GRB in the same field [14]. Short GRBs, on the other hand, typically fade too quickly to observe past the initial burst of gamma rays and hard x-rays. Rapid gravitational-wave transient alerts could enable the observation of optical flashes from short GRBs. An optical counterpart would vastly boost confidence in the gravitational-wave detection and provide the tight sky localization necessary to allow determination of the source’s host galaxy, which leads to a redshift measurement. With both redshift and a coincident gravitational-wave observation, we can produce precision measurements of the Hubble constant [16].

To this end, we have the ambition of reporting candidates not minutes *after* the merger, but seconds *before*. By looking for threshold crossings before the gravitational-wave signal leaves the detection band, it is possible to trade some signal to noise ratio (SNR) for latency. Figure 1 shows projected early trigger rates for NS–NS binaries in Advanced LIGO assuming the event rate predictions in [5].

The gravitational-wave community initiated a project to send alerts when potential gravitational-wave transients are observed. In October 2010, LIGO completed its sixth science run (S6) and Virgo completed its third science run (VSR3). While both LIGO detectors and Virgo were operating, several all-sky detection pipelines operated in a low-latency configuration, namely MBTA, ihope, Coherent WaveBurst, and Omega [20, 21]. The S6 analyses achieved latencies of 30–60 minutes, which were dominated by a human vetting process. Candidates were sent for electromagnetic followup to several telescopes; Swift, ROTSE, TAROT, and Zadko [22, 20] took images of likely sky locations. MBTA achieved the best gravitational-wave trigger generation latencies of 2–5 minutes. We assume that in the advanced detector era the vetting process will be automated, so current trigger generation and telescope actuation would then dominate latency.

Predictive detection of CBCs will require striking a balance between latency and throughput. CBC searches consist of banks of matched filters, or cross-correlations between the data stream and a bank of nominal “template” signals. There are many different implementations of matched filters, but most have high throughput at the cost of high latency, or low latency at the cost of low throughput. The former are epitomized by the overlap-save algorithm for FFT convolution, currently the preferred method in gravitational wave searches. The most obvious example of the latter is the time domain (TD) convolution, which has no latency at all. However, its computational complexity is quadratic in the length of the templates, so it is prohibitively expensive for long templates.

Fortunately, the morphology of inspiral signals can be exploited to offset some of the computational complexity of low latency algorithms. First, the signals evolve slowly in frequency, so that they can be broken into contiguous bandlimited time

*I think we should say also something about how we can also get the full SNR within moments after the coalescence, it is obvious maybe, but we should be explicit*  
*Citation needed for LOOC-UP*  
*Get references for these low-latency pipelines.*

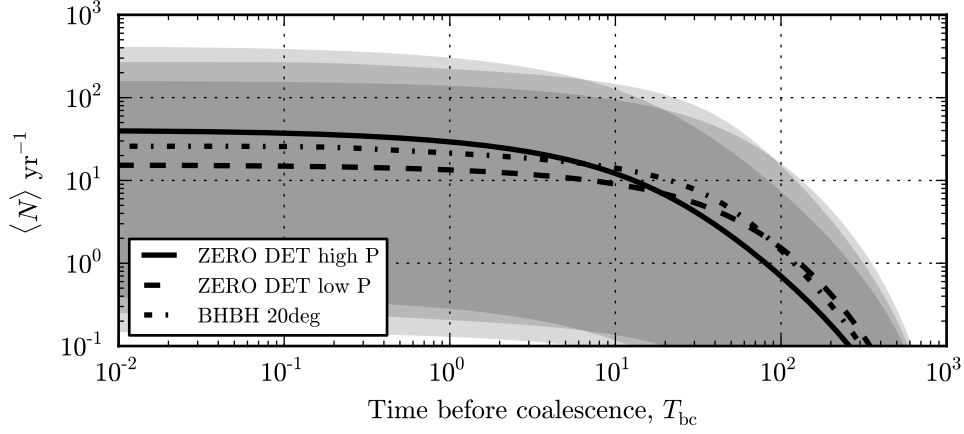


Figure 1: Expected number of NS–NS sources that will be detectable  $T_{\text{bc}}$  seconds before coalescence. The solid line is the most likely yearly rate estimate  $N_{\text{re}}$  for Advanced LIGO and the shaded region is the interval  $N_{\text{low}}$  to  $N_{\text{high}}$  from [5]. Note that assuming SNR 8 is sufficient for detection and that we observe  $N_{\text{re}} = 40$  events per year with a detector having the ZERO\_DET\_high\_P noise model described in [17],  $\sim 10$  sources detectable within 10 seconds of merger and  $\sim 1$  sources detectable within 100 seconds of merger. Other noise models, ZERO\_DET\_low\_P [18] and BHBH\_20deg [19] provide better results for early detection but at the cost of fewer total events observed above SNR 8.

intervals and processed at possibly lower sample rates. Second, inspiral filter banks consist of highly similar templates, admitting principal component analysis to reduce the number of templates. We described a rank-reduction scheme based on singular value decomposition in [23]. We will use both aspects to demonstrate that a very low latency analysis with predictive detection of compact binary sources is possible with current computing resources. Assuming other aspects of gravitational-wave observation latency can be reduced significantly, this should allow the possibility for prompt alerts to be sent to the astronomical community.

The paper is organized as follows. First we provide an overview of our method for detecting compact binary coalescence signals in an early-warning analysis. We then describe the pipeline we have constructed that implements our method. To validate the approach we present results of simulations and conclude with some remarks on what remains to prepare for the Advanced detector era.

## 2. Early warning searches for compact binary coalescence

In this section we describe a decomposition of the compact binary parameter space that reduces low-latency filtering cost sufficiently to allow for the possibility of early-warning detection with modest computing requirements. We expand on the ideas of [24, 25] that describe a multiband decomposition of the compact binary parameter space that resulted in search with  $\sim$ minutes latency in LIGO’s S6 and Virgo’s VSR2 science runs. We combine this with the orthogonal decomposition described in [23] that exploits the redundancy of the template banks.

### 2.1. Conventional CBC matched filter searches

Inspiral signals are parameterized by a set of intrinsic parameters  $\bar{\theta}$  that determine the amplitude and phase evolution of the gravitational wave strain. For systems where the effects of spin can be ignored, the intrinsic parameters are the component masses of the binary,  $\bar{\theta} = (m_1, m_2)$ . Searches for gravitational waves from compact binary coalescence typically employ matched filter banks that discretely sample the possible intrinsic parameters [26]. The filters for the waveforms  $h(t, \bar{\theta})$ , known as templates, are the waveforms weighted by the inverse detector noise amplitude spectral density. To construct a template bank, templates are chosen with discrete signal parameters  $\theta_1, \theta_2, \dots, \theta_N$  to assure a bounded loss of SNR [27, 28]. That is, any possible signal within the search space will have a cross-correlation of  $\geq 0.97$  with at least one template. Such a template bank is said to have a *minimum match* of 0.97. Data are filtered against each template to produce an SNR time-series. Local peak-finding across time and templates determines detection candidates. <sup>‡</sup>

We will denote the  $i^{\text{th}}$  filter with parameters  $\bar{\theta}_i$  as a function of time  $x_i(\tau)$ . In this work we will discuss transformations to a set of filters  $\{x_i(\tau)\}$ . Some of these transformations are not useful or practical over the entire parameter space. For that reason we assume from here onward that the set of filters  $\{x_i(\tau)\}$  refer to a set of near-neighbor filters that can be chosen as a subset of the full parameter space. Several such subsets can be chosen until all of the filters are a member of one local set.

Filtering the detector data  $h(t)$  involves a convolution of the data with the filter. For a unit-normalized filter, and whitened detector data, the result can be interpreted as the signal-to-noise ratio,  $\rho_i(t)$  and is defined as

$$\rho_i(t) = \int x_i(\tau) h(t - \tau) d\tau \quad (1)$$

$$= \int \tilde{x}_i(f) \tilde{h}(f) e^{-2\pi i f t} df, \quad (2)$$

where the second line is a result of the convolution theorem and  $\tilde{x}_i(f)$  is the Fourier transform of  $x_i(\tau)$  as is  $\tilde{h}(f)$  the Fourier transform of  $h(t)$ .

The evaluation of the integrals in (1) and (2) are implemented as sums over sample points for the digitized gravitational-wave detector output. Discrete Fourier transforms can be computed efficiently numerically. For that reason (2) is typically far faster, computationally. To evaluate (1) requires  $\mathcal{O}[N_{x_i} N]$  floating point operations per filter  $x_i$ , where  $N_{x_i}$  is the number of sample points in the filter  $x_i$  and  $N$  is the number of sample points in the data  $h$ . Assuming  $N_t$  filters are required, the total cost is  $\mathcal{O}[N_t N_{x_i} N]$ . However, (2) requires only  $\mathcal{O}[N \log N]$  operations per filter assuming transform lengths that are longer than the filter (i.e.  $N > N_{x_i}$ ), resulting in a total cost of  $\mathcal{O}[N_t N \log N]$ . In most cases  $N_{x_i} \gg \log N$  and the computational savings by choosing the frequency-domain integral form (2) is clear. However, to take full advantage of the computational efficiency of (2) requires an acausal knowledge of the detector data  $h(t)$ , which implies an inherent latency. In contrast, (1) can be updated every time a new sample point of detector data is taken.

<sup>‡</sup> There are two gravitational-wave polarizations,  $+$  and  $\times$ . A given detector will observe a combination of these polarizations that will largely be degenerate with an overall unknown constant phase. This can be maximized over by filtering for quadrature phases and taking the magnitude of the result. For simplicity we will ignore that aspect in this work, as it is straightforward to generalize but not necessary for understanding any of the points that will be made.

*Change this to a harpoon or vector symbol.*

*nvf: I don't think that the splitting of the bank is a crucial point here. Perhaps later if we discuss the up-front cost of the SVD.*

## 2.2. Proposed method

In order to minimize latency we propose using the time-domain convolution presented in (1). However, because the brute-force evaluation of (1) is far too costly to be useful, we will consider an approximation to (1) that can reduce substantially the cost of real-time filtering. This approximation has the form

$$\rho_i(t) \approx \sum_k^{N_{ts}} \sum_j^{N_{uj,k}} \int_{\tau_k}^{\tau_{k+1}} v_{ijk} \sigma_{jk} u_{jk}(\tau) h(t - \tau) d\tau \quad (3)$$

where  $u_{jk}(\tau)$  is an orthogonal basis set of filters spanning the space of  $\{x_i(\tau)\}$  and  $\sigma_{jk} v_{ijk}$  is a tensor relating the filters  $u_{jk}(\tau)$  to the original filter set  $\{x_i(\tau)\}$ . We claim that with a suitable choice of filters  $u_{jk}(\tau)$  one can reduce the computational cost of (1) sufficiently to feasibly search for gravitational waves from compact binary coalescence in real-time. This requires 1) exploiting the time-frequency characteristics of the binary waveforms and 2) exploiting the redundancy of the template bank. We describe our procedure for producing the decomposition in (3) in the remainder of this section.

### 2.2.1. Selectively reducing the sample rate of the data and template waveforms

The first step of the orthogonal decomposition described in (3) is to divide the templates into *time slices*. This is a time-domain analogue to the frequency-domain decomposition described in [24, 25, 29, 30]. A matched filter is constructed for each time slice. The outputs form an ensemble of partial SNR streams. By linearity, these partial SNR streams can be suitably time delayed and summed to reproduce the SNR of the full template. We will show in the next section that this, combined with the singular value decomposition, is sufficient to enable a computationally efficient time-domain search and furthermore is an essential part of an early-warning detection scheme.

For concreteness and simplicity, we will consider an inspiral waveform in the quadrupole approximation, for which the time-frequency relation is

$$f = \frac{1}{\pi \mathcal{M}} \left[ \frac{5}{256} \frac{\mathcal{M}}{-t} \right]^{3/8}. \quad (4)$$

Here,  $\mathcal{M}$  is the chirp mass of the binary in units of time (where  $GM_\odot/c^3 \approx 5\mu\text{s}$ ) and  $t$  is the time relative to the coalescence of the binary [26, 31]. Usually the template is truncated at some prescribed time  $t_0$ , or equivalently frequency  $f_{\text{max}}$ . This is often chosen to correspond to the ISCO defined previously. An inspiral signal will enter the detection band at a low frequency,  $f = f_{\text{low}}$  corresponding to a time  $t_{\text{low}}$ . The template is assumed to be zero outside the interval  $[t_{\text{low}}, t_0]$  and is said to have a duration of  $t_0 - t_{\text{low}}$ . It is critically sampled at a rate of  $2f_{\text{max}}$ .

The monotonic time-frequency relationship of (4) allows us to choose time-slice boundaries that require substantially less bandwidth at early times in the inspiral. Our goal is to reduce the filtering cost of a large fraction of the waveform by computing part of the filter integral at a lower sample rate. Specifically we consider here time slice boundaries with the next highest power-of-two sample rate that critically samples the time sliced filter. The time slices for this template consist of the  $k$  intervals  $(t_k, t_{k-1}]$ ,  $\dots$ ,  $(t_2, t_1]$ ,  $(t_1, t_0]$  sampled at frequencies  $f_{k-1}, \dots, f_1, f_0$  where  $f_0 \geq 2f_{\text{ISCO}}$ ,  $t_0 = t_{\text{ISCO}}$ ,  $f_{k-1} \geq 2f_{\text{low}}$ , and  $t_k \leq t_{\text{low}}$ . An example time-slice design satisfying these constraints for a  $1.4 - 1.4 M_\odot$  binary is shown in table 1.

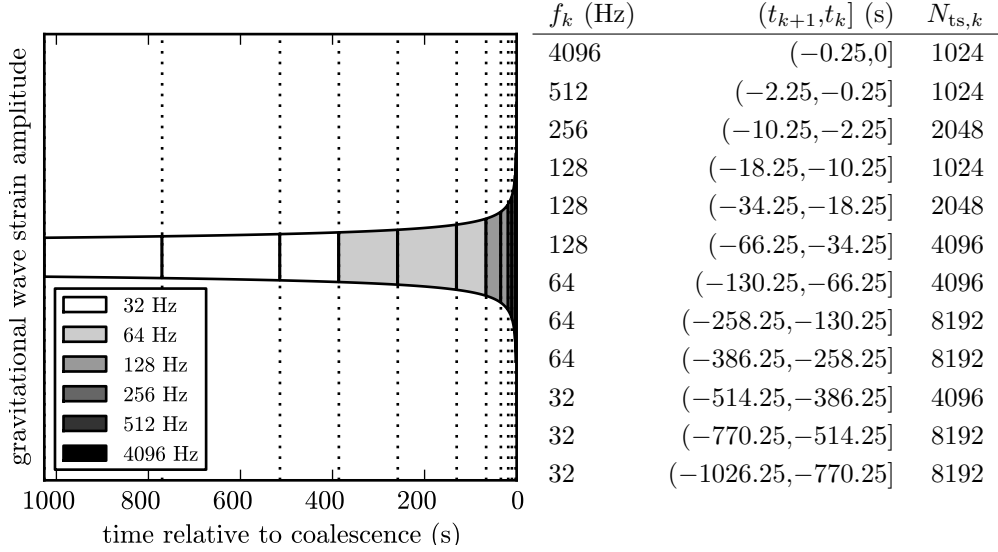


Table 1: Example of nearly critically sampled, power-of-two time slices for a  $1.4 - 1.4 M_\odot$  template extending from  $f_{\text{low}} = 10$  Hz to  $f_{\text{ISCO}} = 1571$  Hz with a time frequency structure given by (4).  $f_k$  is the sample rate of the time slice,  $(t_{k+1}, t_k]$  are the boundaries in seconds preceeding coalescence and  $N_{ts,k}$  are the number of sample points in the  $k^{\text{th}}$  filter.

Rather than applying a unique time-slice decomposition to each template waveform, we find a decomposition that is adequate for the entire set  $\{x_i\}$ . This is facilitated by choosing templates in the set  $\{x_i\}$  with similar chirp masses,  $\mathcal{M}$ . The time-slice decompositions of the filters  $x_i(\tau)$  lead to new filters  $y_{ik}(\tau)$  satisfying

$$x_i(\tau) \approx \sum_k^{N_{ts}} y_{ik}(\tau) \quad (5)$$

where  $N_{ts}$  is the number of time slices required. Note that we write the relationship as approximate since the resampling implementation may have some inherent loss in quality. We note that by construction these filters are orthogonal over the index  $k$  since they are disjoint in time. In the next section we examine how to reduce the number of filters  $y_{ik}(\tau)$  via singular value decomposition to construct a set of filters that is also orthogonal within a given time slice.

**2.2.2. Reducing the number of filters with the singular value decomposition** As described previously, the template banks are, by design, highly correlated. It is possible to greatly reduce the number of filters required to achieve a particular minimum match by designing an appropriate set of orthonormal *basis templates*. A purely numerical technique based on the application of the singular value decomposition (SVD) to inspiral waveforms is demonstrated in [23]. Using the results of [23] we establish that the filters of the previous section can be approximated to

Symbol	Definition
$N$	The number of sample points in the data
$N_{x_i}$	The number of sample points in template $x_i$
$N_t$	The number of templates in the set $\{x_i\}$
$N_{ts}$	The number of time slices
$N_{u_j,k}$	The number of orthogonal filters in time slice $k$
$N_{ts,k}$	The number of samples in time slice $k$
$f_{\max}$	The maximum frequency of a filter
$f_k$	The sample frequency of the $k^{\text{th}}$ filter
$N_r$	The number of sample points in the resample filter

Table 2: Notation used to describe filtering. This table provides a quick reference for symbols used.

high accuracy by the expansion in the singular value basis

$$y_{ik}(\tau) \approx \sum_j^{N_{u_j,k}} \sigma_{jk} v_{ijk} u_j(\tau) \quad (6)$$

where  $\sigma_{jk}$  is the  $j^{\text{th}}$  singular value for the  $k^{\text{th}}$  time slice,  $v_{ijk}$  is an orthogonal matrix for the  $k^{\text{th}}$  time slice and  $u_{jk}$  is a new orthogonal basis filter set for the  $k^{\text{th}}$  time slice. The authors of [23] showed that to high reconstruction accuracy far fewer filters are needed than were in the original template bank. We find that when combined with the time-slice decomposition, the number of SVD filters,  $N_{u_j,k}$  is much smaller than the original number of filters  $N_t$ . We combine (6) with (5) to arrive at (3). In the next section we compute the expected computational cost scaling of this decomposition and compare it with the brute-force implementation in (1) and higher latency FFT methods.

### 2.3. Comparison of computational costs

We now examine the computational cost scaling of the approximate implementation of (1) as (3). An actual implementation of this decomposition in a working analysis pipeline is discussed in the next section along with measured computational requirements. For convenience, table 2 is a recap of the meaning of various symbols used in this calculation.

In table 3 we present the computational cost scaling in floating point operations per sample for common tasks in the pipeline.

The filter bank can be implemented using finite impulse response (FIR) filters, which are just sliding window dot products. If there are  $N_t$  templates of length  $N_{x_i}$  and the data stream contains  $N$  samples, then applying the filter bank requires  $2N_t N_{x_i} N$  operations.

More commonly, the matched filters are implemented using the FFT convolution. This entails applying FFTs to blocks of  $D$  samples, with  $N_{x_i} \leq D$ , each block overlapping the previous one by  $D - N_{x_i}$  samples. There are  $N/(D - N_{x_i})$  such blocks required to filter  $N$  samples of data. Modern implementations of the Cooley-Tukey FFT, such as the ubiquitous `fftw`, require about  $4N \lg N$  operations to evaluate a DFT

Process	ops/sample
FIR matched filter, $N_t$ templates of length $N_{x_i}$	$2N_tN_{x_i}$
FFT matched filter, $N_t$ templates of length $N_{x_i}$ blocks of length $D$	$\frac{4(N_t+1) \lg D + 2N_t}{1 - N_{x_i}/D}$
FIR resampling filter, length $N_r$ for each of $N_t$ templates and sample rates $f_1 < f_2$	$2N_tN_r f_1/f_2$
multiply $M \times L$ real matrix by $L \times 1$ real vector	$2ML$

Table 3: Number of floating point operations per sample (multiplications and divisions) required for a selection of signal processing operations used in LLOID.

of size  $N$  [32]. A  $D$  sample cross-correlation consists of a forward FFT, a  $D$  sample dot product, and an inverse FFT totaling  $8D \lg D + 2D$  operations per block. Per sample, this is  $(8 \lg D + 2)/(1 - N_{x_i}/D)$  operations. As this expression indicates the number of operations increases as the block size  $D$  approaches the filter length  $N_{x_i}$ .

The FIR filter implementation has the advantage that it has no intrinsic latency, whereas the FFT convolution has latency of  $D - N_{x_i}$ . However, the FIR filter implementation has the disadvantage of much greater overhead per sample than the FFT convolution. For a 1ks template sampled at 4096 Hz, the FIR implementation requires about  $N_{x_i}/8 \lg 2N_{x_i} = 2.2 \times 10^4$  times more operations per sample than the FFT implementation. We will now consider the computational cost of the FIR filter implementation described in the previous sections.

It is convenient to express the computational cost of the entire filtering procedure in floating point operations per second flop/s. The cost will be the sum of the cost of the FIR filtering for the orthogonal filter in each time slice plus the cost of reconstructing the original waveforms with matrix operations and resampling. Using the formulas in table 3 we arrive at

$$\text{C.C.} \propto 2 \sum_k^{N_{ts}} (N_{ts,k} + N_t) N_{u_j,k} f_k r \sum_{k, f_k \neq f_{\max}}^{N_{ts}} N_t N_r f_k \quad (7)$$

where the first sum is the total cost for filtering and reconstructing the orthogonal filters  $u_{jk}$ . The second sum is the cost of resampling the reconstructed time-slice outputs to the original sample rate  $f_{\max}$ . It assumes a FIR filter with  $N_r$  sample points. Note that the cost for resampling only occurs for the downsampled time slices. The resampling cost largely depends on  $N_r$ . The computational cost of (7) is dominated by the highest frequency terms in the sum. Comparing just the highest frequency term of (7) with (1) shows  $\mathcal{O}[(N_{ts,\max} + N_t) N_{u_j,\max} f_{\max}]$  flop/s versus  $\mathcal{O}[N_t N_{x_i} f_{\max}]$  flop/s. The full calculation for a particular patch of parameter space is shown in table ??.

FIXME we need to actually say what went into this

### 3. Implementation

In this section we describe an implementation of the method described in section ?? suitable for rapid gravitational-wave searches for compact binary coalescence. The

*This is more commonly known as “overlap-save”. We should find someone else’s operation count and cite it. Drew: Why don’t we change this to an overlap of  $m$  samples so we can see what happens as we increase the overlap to reduce latency.*

*From here down introduce an example, but don’t go through the FFT methods in lloid focus just on the comparison of the TD method without lloid tricks, the FD method without lloid tricks, and the TD lloid method (no comp detection statistic though), remember that the goal is near realtime detection not computational cost*



operations/sample	latency (s.)	method
3,714,580,480	$2.4 \times 10^{-4}$	conventional FIR method (1)
49,937	$1.8 \times 10^3$	conventional FFT method (2)
120,000	$9.0 \times 10^2$	conventional FFT method (2)
120,673	$2.4 \times 10^{-4}$	FIR method with time slices and SVD (3)

Table 4: **FIXME: Operation counts per sample for four different detection methods.**

previous method requires several computations that can be done before the analysis is actually underway. Thus we divide the procedure into two stages 1) an offline planning stage and 2) an online, low latency filtering stage. The offline stage can be done before the analysis is started and updated asynchronously, whereas the online stage must keep up with the detector output and produce search results as rapidly as possible. In the next two subsections we describe what these stages entail.

### 3.1. Planning stage

The choice of filter waveforms and singular value decomposition can be done in advance and will be valid as long as the detector noise spectrum remains roughly constant. New filter waveforms can be computed asynchronously using updated spectrum estimates as they are available.

The planning stage begins with choosing templates that cover the space of mass parameters with a hexagonal grid [33] in order to satisfy the minimum match criterion, which assures a user specified maximum loss in SNR for signals that fall in-between the chosen templates. Next, the templates are subdivided into groups of neighbors called “sub-banks” that are appropriately sized so that each bank can be efficiently handled by a single computer. Dividing the mass space into smaller sub-banks also aids in the computational cost of the singular value decomposition. Using our understanding of the time-frequency evolution of the templates, we choose time slice boundaries as in (??) such that all of the templates within a sub-bank are sub-critically sampled at progressively lower sample rates. Next, the templates within the sub-bank are realized as FIR filter coefficients. For each time slice, the templates are downsampled to the appropriate sample rate. Finally, the SVD is applied to each time slice in order to produce a set of orthogonal FIR filters and a reconstruction matrix that maps them back to the original templates. The downsampled orthogonal FIR filter coefficients, the reconstruction matrix, and the time slice boundaries are all saved to disk.

### 3.2. Filtering stage

We have implemented a prototype of the low latency filtering stage using an open source signal processing environment called GStreamer [34]. GStreamer is a vital component of many Linux systems, providing media playback, authoring, and streaming on devices from cell phones to desktop computers to streaming media servers. It turns out that it is also uniquely suited to gravitational wave data analysis. In our application, GStreamer excels at queueing, synchronizing, adding, and bookkeeping many different signals at different sample rates. It also provides some useful signal processing primitives such as decimators, FIR filters, and interpolators.

Most importantly, it permits us to recruit all of the host system’s CPUs without having to write complicated and error-prone multithreaded code.

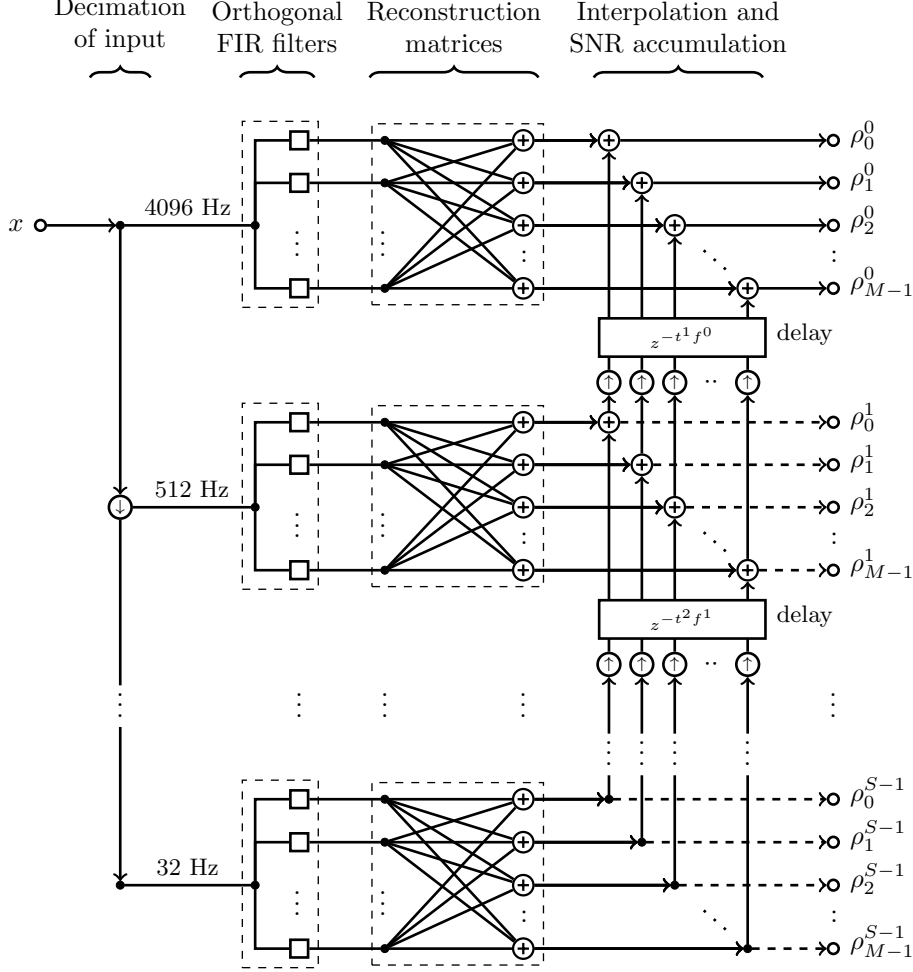


Figure 2: Schematic of LLOID pipeline illustrating signal flow. Circles with arrows represent upsampling  $\uparrow$  or downsampling  $\downarrow$ . Circles with plus signs represent summing junctions  $\oplus$ . Squares  $\square$  stand for FIR filters. Sample rate decreases from the top of the diagram to the bottom.

The filter pipeline consists of six distinct stages.

**3.2.1. Decimation** First, the sample rate of the whitened detector data is reduced to successively lower sample rates by decimation. Decimation involves applying an antialiasing filter to the data, and then downsampling by deleting samples. We use a 192-tap FIR decimator provided by GStreamer.

**3.2.2. Time delays** Each decimated detector data stream becomes the input for one time slice, but it must be appropriately delayed.

### 3.2.3. Orthogonal FIR filters

### 3.2.4. Reconstruction

### 3.2.5. Interpolation

### 3.2.6. SNR accumulation

## 4. Results

[ I would like to see this section changed to emphasize the low-latency time domain filtering 1. Use a live simulated white noise source (this ignores the latency of whitening, but that goes beyond the scope of this paper, and we mention this and perhaps suggest some exaporation 2. Use TD filtering of N templates 3. Present the performance and latency and provide estimates for number of cores required for realtime ALIGO search based on current infrastructure 4. Compute the impulse responses of the templates and histogram the SNR loss for the SVD and time/slice/resampling 5. Put a tee after some of the low frequency stages and perhaps test the possibility of predictive filtering (also with latency measurments) This will remove the complications of running a full pipeline. We need to sort things out before we do that, and this paper should not be delayed any further. The above tests will make the point we need to make ]

### 4.1. Detector noise characteristics

We tested the new detection method with mock Advanced LIGO data having a power spectrum prescribed by the “zero detuning, high power” noise model in [35].

## 5. Conclusions

We have demonstrated a computationally feasible procedure for the rapid detection of gravitational waves emitted during the coalescence of neutron stars and stellar-mass black holes. These sources are expected to produce prompt electromagnetic signals and may be the progenitors of some short hard gamma-ray bursts. Rapid alerts to the broader astronomical community may improve the chances of detecting an electromagnetic counterpart in bands from X-ray down to radio. We antipate requiring no more than 1000 modern computer cores to analyze a four-detector network of gravitational-wave data for such systems assuming spin effects can be ignored in parameterizing the expected signals [?]. This is within the current computing capabilities of the LSC Data Grid [36].

The algorithm we described has no intrinsic latency. However, there are fundamental and practical latencies associated with the analysis and detection procedure. For example, the LIGO detectors, data acquisition is synchronized to a 1/16 Hz cadence introducing an up-front latency of 125 ms. Data aggregation from the observatories will travel over various networks, each capable of high bandwidth but perhaps only modest latency. This could amount to a similar latency of  $\sim 100$  ms. Lastly, unless a realtime infrastructure is adopted post data acquisition, it is likely that there will be an inherent latency introduced by such infrastructure. We have shown a prototype implementation using `gstlal` that is capable of  $\sim 1$  s latency. In

our opinion, significant work would have to be done in order to improve upon this number. However, it should be considered for third generation detector design. For example, a tighter integration of analysis and data acquisition would be beneficial.

We have omitted discussion of source localization. Localization is known to be poor for signals of low SNR [37]. However, one should not immediately dismiss the practical usefulness of a poorly localized source. Even with poor localization, it should be possible to begin downselecting what observatories could view a potential signal and for such observatories to begin any necessary prerequisite activities. In future works we will explore more rigorously the pointing prospects with realistic simulations using the infrastructure and techniques described in this work.

Latency budget, including ‘before’ and ‘after’ quotes for:

- Data acquisition
- Calibration
- Data aggregation
- Analysis
- Localization
- Alert
- Telescope actuation
- Total

Future work:

- Sub-solar mass search
- Hierarchical detection

## Acknowledgments

LIGO was constructed by the California Institute of Technology and Massachusetts Institute of Technology with funding from the National Science Foundation and operates under cooperative agreement PHY-0107417. CH would like to thank Ilya Mandel for many discussions about rate estimates and the prospects of early detection. NF would like to thank Alessandra Corsi for illuminating discussions on astronomical motivations.

This paper has LIGO Document Number LIGO-P0900004-v3.

## References

- [1] Advanced LIGO [www.advancedligo.mit.edu](http://www.advancedligo.mit.edu)
- [2] Advanced Virgo <https://www.cascina.virgo.infn.it/advirgo>
- [3] GEO 600 <http://www.geo600.org>
- [4] Large-scale Cryogenic Gravitational wave Telescope (LCGT) <http://gw.icrr.u-tokyo.ac.jp/lcgt/>
- [5] Abadie J *et al.* 2010 *Class. Quant. Grav.* URL <http://iopscience.iop.org/0264-9381/27/17/173001>
- [6] Shibata M and Taniguchi K 2008 *Phys. Rev. D* **77** 084015–+ (*Preprint* [arXiv:0711.1410](https://arxiv.org/abs/0711.1410))
- [7] Lee W H, Ramirez-Ruiz E and Granot J 2005 *Astrophys. J.* **630** L165–L168
- [8] Nakar E 2007 *Phys. Rept.* **442** 166–236 (*Preprint* [astro-ph/0701748](https://arxiv.org/abs/astro-ph/0701748))
- [9] Sari R and Piran T 1999 *The Astrophysical Journal* **520** 641 URL <http://stacks.iop.org/0004-637X/520/i=2/a=641>



- 2526matchBoolean%253Dtrue%2526rowsPerPage%253D50%2526searchField%253DSearch+All
- [33] Cokelaer T 2007 *Phys. Rev. D* **76** 102004
  - [34] GStreamer: open source multimedia framework URL <http://gstreamer.freedesktop.org>
  - [35] Shoemaker D 2009 Advanced ligo anticipated sensitivity curves Tech. rep. URL <https://dcc.ligo.org/cgi-bin/DocDB/ShowDocument?docid=2974>
  - [36] The LSC Data Grid URL <https://www.lsc-group.phys.uwm.edu/lscdatagrid/>
  - [37] Fairhurst S 2009 *New Journal of Physics* **11** 123006