# Description of Polgraw All-Sky Search GPGPU program
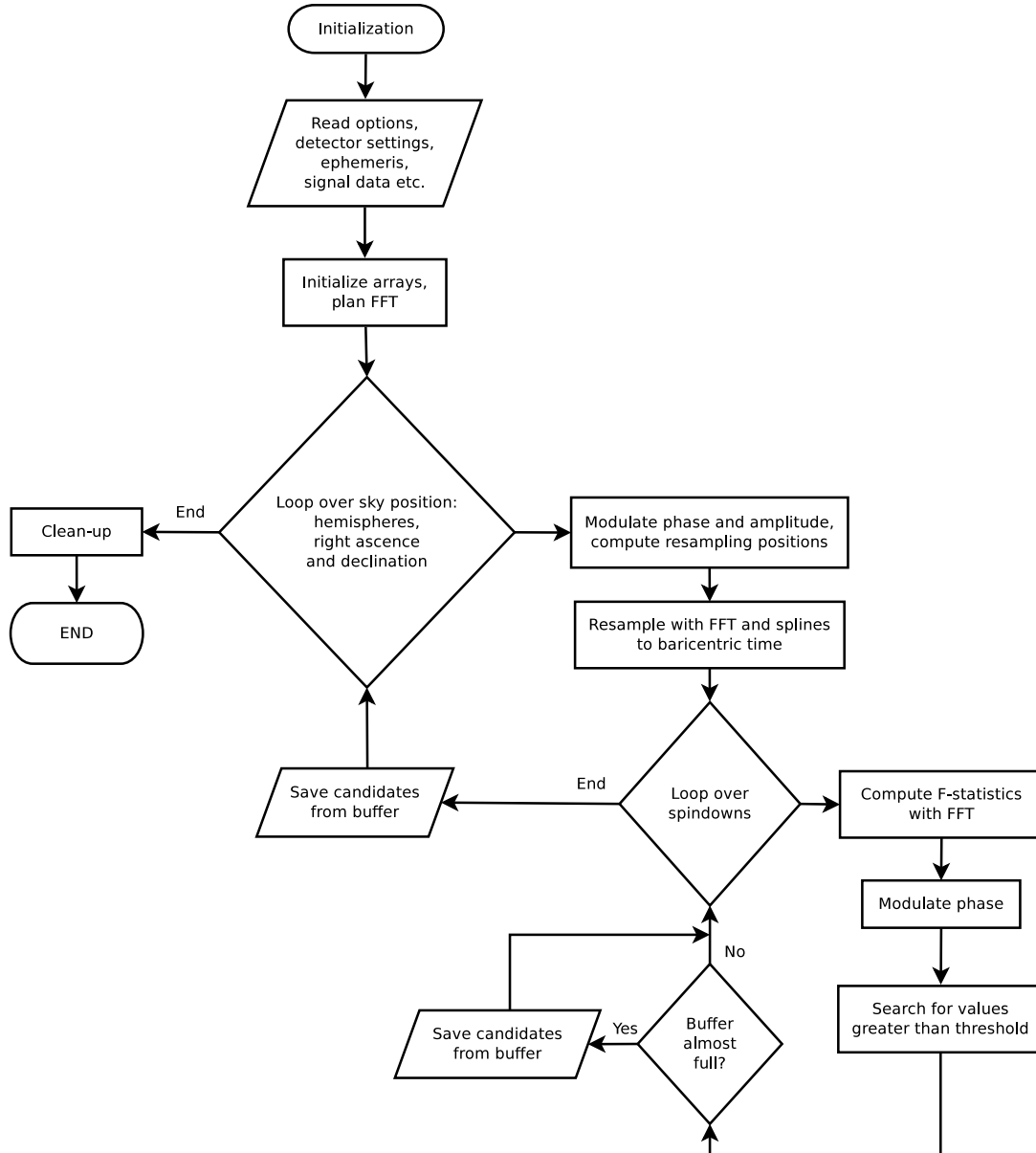
Jan Bolek

# Contents

# 1  Introduction

The purpose of this article is to introduce one into structure, working principle and techniques used in gravitational waves (abbr. GW) All-Sky Search program.

The main purpose of described program is to find candidate GW signals parameters using $\mathcal{F}$-statistic. The parameters consist of frequency, spin-down (first derivative of frequency), positions on sky (e.g. right ascension and declination) and signal-to-noise ratio (SNR). Further theoretical explanation can be found in [1].

## 1.1  The algorithm

Figure 1: Flow-chart of the algorithm

The searching algorithm consists of iterating over sky positions (i.e. hemispheres, right ascension and declination transformed into integral-valued grid), frequency and one spin-down parameter in order to modulate signal with appropriate function using this parameters and evaluate so-called $\mathcal{F}$-statistic. The trick is that $\mathcal{F}$-statistic may be computed using Fast Fourier Transform (FFT), thus iterating over frequency loop can be avoided, which results in great performance improvement.

In order to apply above method detector signal must be defined over baricentric time (i.e. time with respect to solar system baricenter), so there must be performed interpolation to precomputed baricentric time positions. Those positions are known from detector's and Earth's ephemeris.

Because of usage of large vectors ($\sim 2^{19}$ elements) and Fourier transform the program was written for GPUs using CUDA architecture. Thanks to use of massive parallelization there was achieved speed-up over CPU version of about 50-100 times, depending on GPU type.

# 2 Code explanation

## 2.1 Data structures

For convenience arrays and other variables, e.g. settings and detector constants are stored in structures defined in file `struct.h`. Basic structures are:

- `struct _comm_line_opts` (typedef `Command_line_opts`)
  options read from command-line

- `struct _detector_settings` (typedef `Detector_settings`)
  parameters and settings for detector and for computations (array lengths etc.)

- `struct _search_range` (typedef `Search_range`)
  range of search

- `struct _fft_plans` (typedef `FFT_plans`)
  plans (handles) for cuFFT

- `struct _ampl_mod_coeff` (typedef `Ampl_mod_coeff`)
  amplitude modulation function coefficients

## 2.2 "main" function

`main()` function consists of other functions calls – mostly initialization. The most important part is:

Listing 1: main.c

```
1   /*
2   ############ Main job ###############
```

```
3       */
4
5       search(
6              &sett,
7              &opts,
8              &s_range,
9              &arr,
10             &fft_plans,
11             &amod,
12             &Fnum,
13             cu_F
14           );
```

i.e. the worker function call.

## 2.3    Initialization

Initialization consists of the following steps:

1. handling command line options (`handle_opts()`)

2. checking if output directory exists (and creating it if needed)

3. reading grid generating matrix (`read_grid()`)

4. reading detector settings and physical constants (`settings()` from `settings.c`)

5. initializing amplitude modulation function coefficients for Virgo detector (`rogcvir()`)

6. allocating arrays and reading input arrays (`init_arrays()`)

7. determining searching range (`set_search_range()`)

8. planning FFT (`plan_fft()`)

9. reading checkpoints (for resuming execution at given point) (`read_checkpoints()`)

All those functions, except for `settings()`, are defined in `init.c`.

## 2.4    Main job

Function `search()` is defined in file `jobcore.cu`. It contains main loop over sky positions and jobcore() function call.

jobcore() function performs computations outside of spin-down loop and in the spin-down loop itself.

## 2.5 Sky loop internals

### 2.5.1 Preparing variables

First, we compute real angular coordinates from linear grid coordinates:

Listing 2: jobcore.cu

```
202     //change  linear  (grid)  coordinates  to  real  coordinates
203     lin2ast (al1/sett->oms, al2/sett->oms, pm, sett->sepsm, sett->cepsm,
204             &sinalt, &cosalt, &sindelt, &cosdelt);
205
206     // calculate  declination  and  right  ascencion
207     // which  will  be  written  in  file  as  candidate  signal  sky  positions
208     sgnlt[2] = asin (sindelt);
209     sgnlt[3] = fmod (atan2 (sinalt, cosalt)+2.*M_PI, 2.*M_PI);
```

Then we compute amplitude modulation functions $a(t)$ and $b(t)$ (`arr->cu_aa` and `arr->cu_bb`, respectively) using GPU. In this step we need to normalize those function by square root of their sum of squares (i.e. $a_i = a_i' / \sqrt{\sum_t a_i'^2}$ ), which is performed using parallel reduction in kernels `reduction_sumsq()` and `reduction_sum()`.

Listing 3: jobcore.cu

```
212     // amplitude  modulation  function
213     modvir_gpu(sinalt, cosalt, sindelt, cosdelt, sett->sphir,
214               sett->cphir, arr->cu_aa, arr->cu_bb, sett->N, arr);
```

Then kernel `shift_time_mod()` computes shifted (baricentric) time, modulates signal with $a(t)$, $b(t)$ and phase factor (thus obtaining two vectors, `xa` and `xb`), computes shifted time position `tshift[i])` for spline interpolation and pads zeros to nearest power of two (`nfft`).

Listing 4: jobcore.cu

```
216     //calculate  detector  positions  with  respect  to  baricenter
217     nSource[0] = cosalt*cosdelt;
218     nSource[1] = sinalt*cosdelt;
219     nSource[2] = sindelt;
220
221     shft1 = 0.;
222     for (j=0; j<3; j++)
223        shft1 += nSource[j]*arr->DetSSB[j];
224     het0 = fmod (nn*sett->M[8]+mm*sett->M[12], sett->M[0]);
225
226     //compute  shifted  time,  modulate  phase  and  pad  zeros  to  size  'nfft'
227     shift_time_mod<<<(sett->nfft + BLOCK_SIZE - 1)/BLOCK_SIZE, BLOCK_SIZE
           >>>(shft1, het0, nSource[0], nSource[1], nSource[2],
228        arr->cu_xDat, arr->cu_xa, arr->cu_xb, arr->cu_shft,
229        arr->cu_shftf, arr->cu_tshift, arr->cu_aa, arr->cu_bb, arr->
               cu_DetSSB,
230        sett->oms, sett->N, sett->nfft, sett->interpftpad);
```

Listing 5: kernels.cu

```
42  __global__ void shift_time_mod(...) {
43     int i = blockIdx.x * blockDim.x + threadIdx.x;
```

```
44      if (i < N) {
45          double S = ns0 * DetSSB[i*3+0] + ns1 * DetSSB[i*3+1] + ns2 *
                DetSSB[i*3+2];
46          shft[i] = S;
47          shftf[i]= S - shft1;
48
49          /* phase mod */
50          double phase = -het0*i - oms * S;
51          double c = cos(phase), s = sin(phase);
52          xa[i].x = xDat[i] * aa[i] * c;
53          xa[i].y = xDat[i] * aa[i] * s;
54          xb[i].x = xDat[i] * bb[i] * c;
55          xb[i].y = xDat[i] * bb[i] * s;
56
57          //calculate time positions for spline interpolation
58          tshift[i] = interpftpad * ( i - shftf[i] );
59      } else if (i < nfft) {
60          xa[i].x = xa[i].y = xb[i].x = xb[i].y = 0;
61      }
62  }
```

### 2.5.2   Resampling and spline interpolation to baricentric time

First we compute forward FFT of length `nfft` (with zeros padded from original length `N`). Then we pad Fourier transform with zeros to size 2*`nfft` (remembering that in unshifted FFT positive frequencies are placed from 0th bin to n/2-1, and negative are placed from n/2 to n-1, so we basically insert zeros "in-between") and take Fourier transform back to obtain signal samples at twice the original sampling frequency (see figure 2.) Then the vector is scaled because of property of FFT (i.e. elements after two FFTs are multiplied by size of vector).
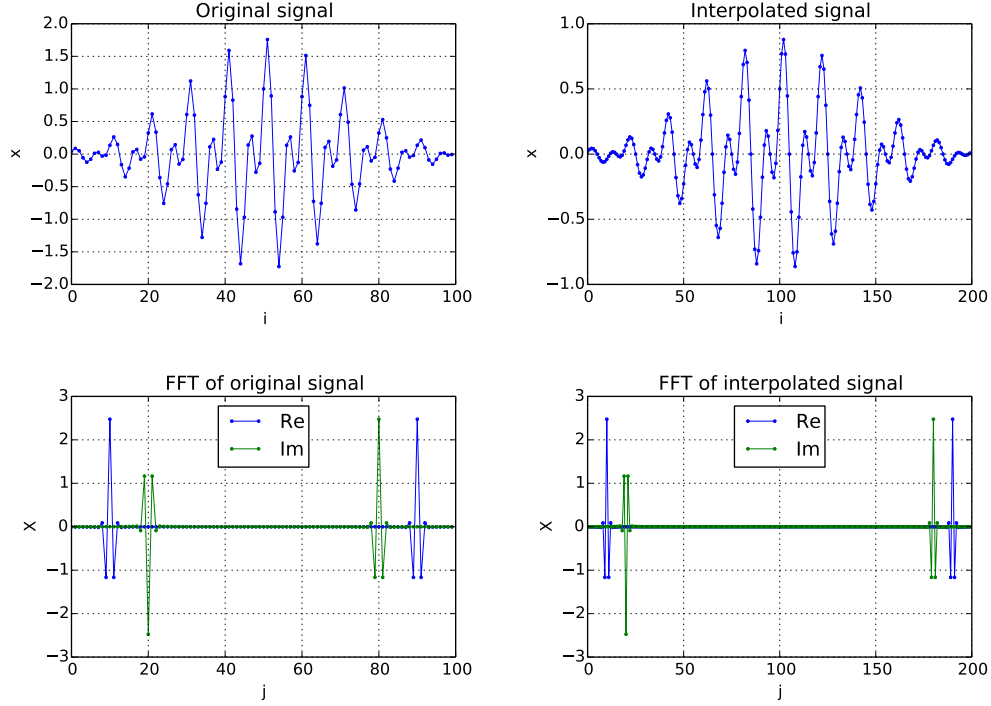
Listing 6: jobcore.cu

```
240     int nyqst = (sett->nfft+2)/2; // Nyquist frequency
241
242     // fft
243     cufftExecZ2Z(plans->pl_int, arr->cu_xa, arr->cu_xa, CUFFT_FORWARD);
244     cufftExecZ2Z(plans->pl_int, arr->cu_xb, arr->cu_xb, CUFFT_FORWARD);
245
246     //shift frequencies and remove those over Nyquist
247     resample_postfft<<<(sett->Ninterp + BLOCK_SIZE - 1)/BLOCK_SIZE,
            BLOCK_SIZE>>>(arr->cu_xa, arr->cu_xb, sett->nfft, sett->Ninterp,
            nyqst);
248     CudaCheckError();
249
250     // fft back
251     cufftExecZ2Z(plans->pl_inv, arr->cu_xa, arr->cu_xa, CUFFT_INVERSE);
252     cufftExecZ2Z(plans->pl_inv, arr->cu_xb, arr->cu_xb, CUFFT_INVERSE);
253
254     //scale fft
255     ft = (double)sett->interpftpad / sett->Ninterp; //scale FFT
256     scale_fft<<<(sett->Ninterp + BLOCK_SIZE - 1)/BLOCK_SIZE, BLOCK_SIZE
            >>>(arr->cu_xa, arr->cu_xb, ft, sett->Ninterp);
257     CudaCheckError();
```

Figure 2: FFT interpolation example



Next step is spline interpolation in which we use cubic splines with natural boundary conditions, meaning that second derivatives at boundaries ($x_0$ and $x_n$) are set to zero. In our case we also have other assumption: $h = x_{i+1} - x_i = 1$.

Spline coefficients are computed by solving linear equation:

$$
\begin{bmatrix}
4 & 1 & & & & & \\
1 & 4 & 1 & & & & \\
 & 1 & \ddots & \ddots & & & \\
 & & \ddots & \ddots & 1 & & \\
 & & & 1 & 4 & 1 & \\
 & & & & 1 & 4
\end{bmatrix}
\begin{bmatrix}
z_1 \\ z_2 \\ z_3 \\ \vdots \\ z_{n-3} \\ z_{n-2} \\ z_{n-1}
\end{bmatrix}
=
\begin{bmatrix}
6(b_1 - b_0) \\
6(b_2 - b_1) \\
6(b_3 - b_2) \\
\vdots \\
6(b_{n-3} - b_{n-4}) \\
6(b_{n-2} - b_{n-3}) \\
6(b_{n-1} - b_{n-2})
\end{bmatrix}
$$

where:

- $n$ is number of intervals in interpolated function, so $n + 1$ is number of given points (knots)

- $z_i$ are unknown coefficients (second derivatives at knots, to be precise),

- $n - 2$ is rank of set of equations (we know from natural boundary conditions that $z_0 = z_n = 0$)

- $b_i$ are differences of function values in consecutive knots: $b_i = (y_{i+1} - y_i)$

7

Matrix of above shape is called *tridiagonal matrix*. We solve this set of equations using `cusparseZgtsv` function included in the `cuSPARSE` package.

Having solved above equations for $z_i$, we can compute interpolated value $y'(x)$ using following formula:

$$y(x) = S_i(x) = \frac{z_{i+1}}{6}(x - x_i)^3 + \frac{z_i}{6}(x_{i+1} - x)^3 + \left(y_{i+1} - \frac{z_{i+1}}{6}\right)(x - x_i) + \left(y_i - \frac{z_i}{6}\right)(x_{i+1} - x)$$

where $i$ is number of spline interval $S_i$ which contains point $x$: $i = \lfloor x \rfloor$.

Function `gpu_interp()` encapsulates above algorithm. It is defined in file `spline_z.cu`.

Listing 7: jobcore.cu

```
261     //interpolation  of  xa
262     gpu_interp ( arr - > cu_xa ,       //input  data
263              sett - > Ninterp ,        //input  data  length
264              arr - > cu_tshift ,       //output  time  domain
265              arr - > cu_xar ,          //output  values
266              sett - > N ,              //output  data  length
267              arr - > cu_d ,            //diagonal
268              arr - > cu_dl ,           //upper  diagonal
269              arr - > cu_du ,           //lower  diagonal
270              arr - > cu_B );           //coefficient  matrix
271
272     //interpolation  of  xb
273     gpu_interp ( arr - > cu_xb ,       //input  data
274              sett - > Ninterp ,        //input  data  length
275              arr - > cu_tshift ,       //output  time  domain
276              arr - > cu_xbr ,          //output  values
277              sett - > N ,              //output  data  length
278              arr - > cu_d ,            //diagonal
279              arr - > cu_dl ,           //upper  diagonal
280              arr - > cu_du ,           //lower  diagonal
281              arr - > cu_B );           //coefficient  matrix
```

## 2.6   Spin-down loop internals

In spin-down loop signal resampled in last step is phase-modulated by factor dependent on current spin-down parameter and shifted baricentric time.

Listing 8: jobcore.cu

```
325         phase_mod_2 <<<( sett - > N + BLOCK_SIZE - 1)/ BLOCK_SIZE , BLOCK_SIZE
                >>>(
326         arr - > cu_xa_f , arr - > cu_xb_f ,
327         arr - > cu_xar_f , arr - > cu_xbr_f ,
328         het1 , sgnlt [1] , arr - > cu_shft ,
329         sett - > N );
```

Next step is the core of whole algorithm: we compute two integrals used in evaluating $\mathcal{F}$-statistic and at the same time perform interpolation to get twice (or more) denser sampling rate. Interpolation can be done with two methods:

1. *interbinning*, which consists in computing values in between the bins by approximation using following formula[1]:
$$x_{i+1/2} = (x_i - x_{i+1})/\sqrt{2}$$

2. *FFT interpolation*, i.e. padding with zeros to size `fftpad`-times larger than original array. This method requires computing twice as long FFT.

Listing 9: jobcore.cu

```
334        if (opts->fftinterp == INT) { //interpolation by interbinning
335            //pad zeros from N to nfft
336            pad_zeros<<<(sett->nfft - sett->N + BLOCK_SIZE - 1)/
                   BLOCK_SIZE, BLOCK_SIZE>>>(
337                arr->cu_xa_f+sett->N, arr->cu_xb_f+sett->N, sett->nfft -
                       sett->N);
338            CudaCheckError();
339
340            // FFT length nfft
341            CUFFT_EXEC_FFT(plans->plan, arr->cu_xa_f, arr->cu_xa_f,
                   CUFFT_FORWARD);
342            CUFFT_EXEC_FFT(plans->plan, arr->cu_xb_f, arr->cu_xb_f,
                   CUFFT_FORWARD);
343
344            // make a gap between every bin
345            interbinning_gap<<<(sett->nfft/2 + BLOCK_SIZE - 1)/
                   BLOCK_SIZE, BLOCK_SIZE>>>(
346                arr->cu_xa_f, arr->cu_xb_f, arr->cu_xa2_f, arr->cu_xb2_f,
                       sett->nfft);
347            CudaCheckError();
348
349            //interpolate values between every bin
350            interbinning_interp<<<( (sett->nfft-2)/2 + BLOCK_SIZE - 1)/
                   BLOCK_SIZE, BLOCK_SIZE>>>(
351                arr->cu_xa2_f, arr->cu_xb2_f, sett->nfft);
352            CudaCheckError();
353
354            cu_xa_final = arr->cu_xa2_f;
355            cu_xb_final = arr->cu_xb2_f;
356
357        } else { //interpolation by FFT (fftpad-times-longer FFT)
358            //pad zeros from N to nfft*fftpad
359            pad_zeros<<<(sett->nfft*sett->fftpad - sett->N + BLOCK_SIZE
                   - 1)/BLOCK_SIZE, BLOCK_SIZE>>>(
360                arr->cu_xa_f+sett->N, arr->cu_xb_f+sett->N, sett->nfft*
                       sett->fftpad - sett->N);
361            CudaCheckError();
362
363            //fft length fftpad*nfft
364            CUFFT_EXEC_FFT(plans->plan, arr->cu_xa_f, arr->cu_xa_f,
                   CUFFT_FORWARD);
365            CUFFT_EXEC_FFT(plans->plan, arr->cu_xb_f, arr->cu_xb_f,
                   CUFFT_FORWARD);
366
367            cu_xa_final = arr->cu_xa_f;
368            cu_xb_final = arr->cu_xb_f;
369        }
```

---

[1] I'm not sure about the sign, but later we use moduli of these values, so it's irrelevant.

Finally, having computed $F_a$ and $F_b$ integrals (as defined in [1], equation (5.5)) we can compute $\mathcal{F}$-statistic ([1], eq. (3.8)) and normalize it.

Listing 10: jobcore.cu

```
374          compute_Fstat <<<( sett ->nmax -sett ->nmin + BLOCK_SIZE - 1)/
                 BLOCK_SIZE , BLOCK_SIZE >>>(
375                               cu_xa_final + sett ->nmin ,
376                               cu_xb_final + sett ->nmin ,
377                               cu_F + sett ->nmin ,
378                               sett ->crf0 ,
379                               sett ->nmax - sett ->nmin );
```

With assumption that noise present in signal is white noise, normalization can be done by simply dividing $\mathcal{F}$-statistic by signal's variance. Without that assumption $\mathcal{F}$-statistic is divided into blocks (size, for example, 4096 elements) and normalized by it's sum in every block[2]:

Listing 11: jobcore.cu

```
383          if ( sett ->sig2 < 0.0) { // when noise is not white-noise
384            FStat_gpu ( cu_F+sett ->nmin , sett ->nmax - sett ->nmin , NAV , arr
                 ->cu_mu , arr ->cu_mu_t );
385          } else { // when noise is white-noise
386            kernel_norm_Fstat_wn <<<( sett ->nmax - sett ->nmin + BLOCK_SIZE
                 - 1)/BLOCK_SIZE , BLOCK_SIZE >>>(
387                               cu_F + sett ->nmin ,
388                               sett ->sig2 ,
389                               sett ->nmax - sett ->nmin );
390          }
```

## 2.7  Finding and saving candidates data

Next step consist in finding $\mathcal{F}$-statistic values that exceed given threshold and save found parameters in output file.

Searching is performed in parallel: every thread checks if value lying under its index is greater than neighboring values. If so it increments buffer counter and saves parameters to local buffer (i.e. spin-down-loop-internal buffer).

Listing 12: jobcore.cu

```
392          //zero candidates count
393          cudaMemset ( arr ->cu_cand_count , 0, sizeof(int));
394
395          /* FINDING AND SAVING CANDIDATES */
396          find_candidates <<<( sett ->nmax - sett ->nmin + BLOCK_SIZE - 1)/
                 BLOCK_SIZE , BLOCK_SIZE >>>(
397            cu_F ,                       // F-statistic array
398            arr ->cu_cand_params , // local parameters buffer
399            arr ->cu_cand_count ,  // local found candidates count
400            opts ->trl ,                 // threshold value
401            sett ->nmin ,                // minimum F-statistic array index
```

---

[2]Important: 'block' length `NAV` described in this paragraph isn't the same as 'block' on the GPU grid length `BLOCK_SIZE`.

```
402              sett->nmax,                // maximum -||-
403              sett->fftpad,              // FFT padding for calculation of freq.
404              sett->nfft,                // FFT length
405              sgnl0,                     // base frequency parameter
406              sett->nd,                  // number of degrees of freedom
407              sgnlt[1],                  // spindown
408              sgnlt[2],                  // sky\
409              sgnlt[3]);                 // positions
```

```
204  // This may be confusing, but three "if"s in kernel below
205  // are faster than one "if" with one, big, combined condition.
206  // Operations in every "if" are the same, so they are wrapped in macro.
207
208  // parameters are:
209  // [frequency, spindown, position1, position2, snr]
210  #define ADD_PARAMS_MACRO \
211        int p = atomicAdd(found, 1); \
212        params[p*NPAR + 0] = 2.0*M_PI*(idx)/((double)fftpad*nfft)+sgnl0; \
213        params[p*NPAR + 1] = sgnl1; \
214        params[p*NPAR + 2] = sgnl2; \
215        params[p*NPAR + 3] = sgnl3; \
216        params[p*NPAR + 4] = sqrt(2*(F[idx]-ndf));
217
218
219  __global__ void find_candidates(...) {
220
221     int idx = blockIdx.x * blockDim.x + threadIdx.x + nmin;
222
223     if (idx > nmin && idx < nmax && F[idx] >= val && F[idx] > F[idx+1] &&
            F[idx] > F[idx-1]) {
224        ADD_PARAMS_MACRO
225     } else if (idx == nmin && F[idx] >= val && F[idx] > F[idx+1]) {
226        ADD_PARAMS_MACRO
227     } else if (idx == nmax-1 && F[idx] >= val && F[idx] > F[idx-1]) {
228        ADD_PARAMS_MACRO
229     }
230  }
```

Finally found candidates parameters are copied to global buffer. If size left in global buffer is less than local buffer size (i.e. next searching could cause overflow) then we save data to output file. Saving is also performed in between first and second sky positions loop as there may not be enough found candidates to flush buffer to disk.

Listing 14: jobcore.cu

```
414         if (cand_count>0) { //if something was found
415            printf ("\nSome signals found in %d %d %d %d\n", pm, mm, nn,
                   ss);
416            cudaMemcpy(arr->cu_cand_buffer + *cand_buffer_count * NPAR,
417                       arr->cu_cand_params,
418                       sizeof(FLOAT_TYPE) * cand_count*NPAR,
419                       cudaMemcpyDeviceToDevice);
420            *cand_buffer_count += cand_count;
421
422            // check if it's time to save results
423            // (place left in buffer is less than one-spindown buffer
                   length)
```

```
424          if (* cand_buffer_count >= arr -> cand_buffer_size - arr ->
                cand_params_size) {
425          save_candidates ( arr -> cu_cand_buffer , arr -> cand_buffer ,
                cand_buffer_count , outname );
426          }
427       }
```

# 3    References

[1]   P. Astone, K. Borkowski, P. Jaranowski, M. Pietka, and A. Królak. "Data analysis
      of gravitational-wave signals from spinning neutron stars. V. A narrow-band all-sky
      search". In: *Physical Review D* 82, 022005 (2010).