Presentation

# Developing Medical imaging application across GPU, FPGA, and CPU using oneAPI

Presenters: Wang Yong (yong4.wang@intel.com), Scott Wang(scott.wang@intel.com)

With contribution from: Zhou Yongfa , Wang Yang, Xu Qing, Wang Chen

intel.

# Table of Content
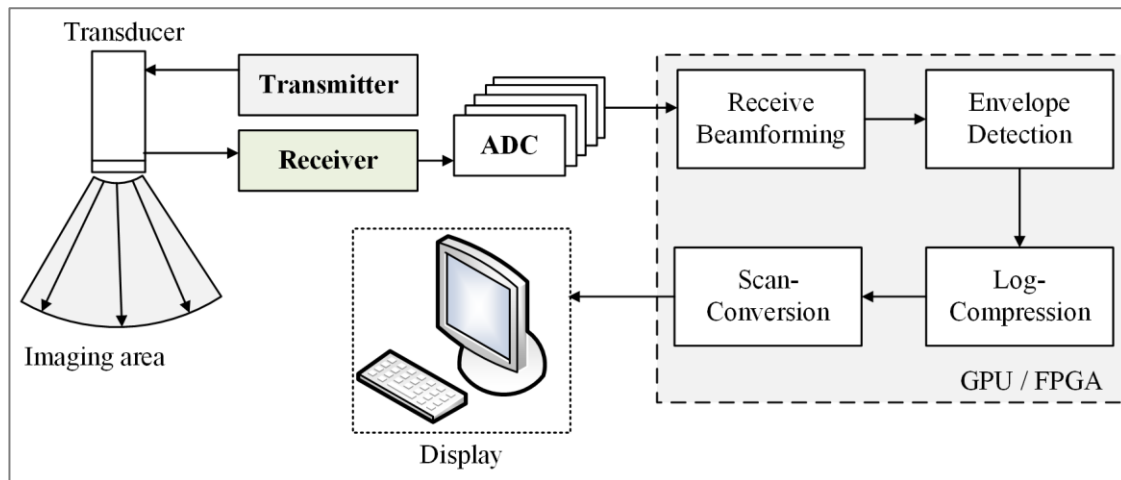
- Background

- Code Migration

- Beamforming Optimization on GPU

- Beamforming implementation on FPGA

- Results and performance

# Background
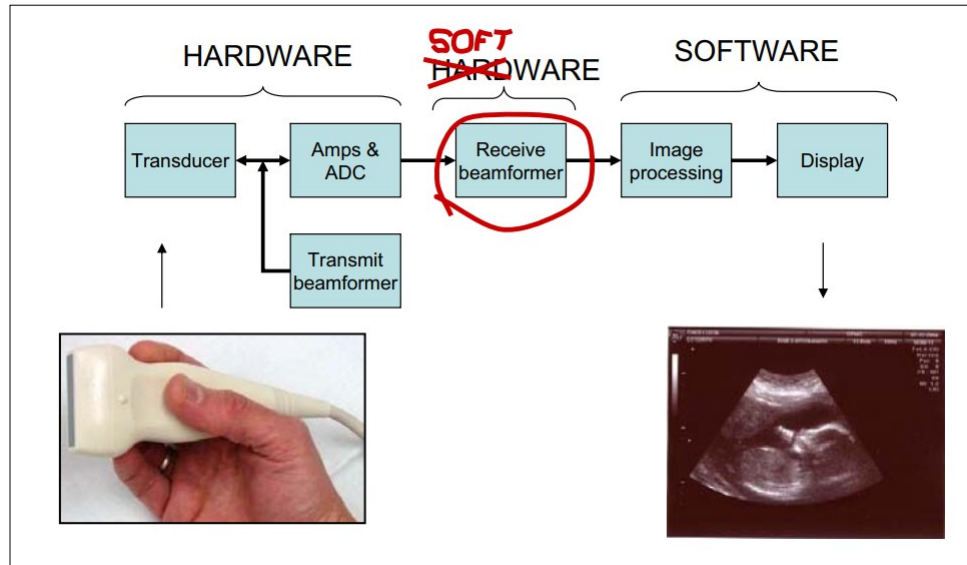
intel.

# What is SUPRA and why we need it?

SUPRA is an open-source pipeline for fully software defined ultrasound processing.

- https://github.com/IFL-CAMP/supra

intel.

# What is software beamforming?

SUPRA contains standard medical ultrasound software beamforming algorithms.
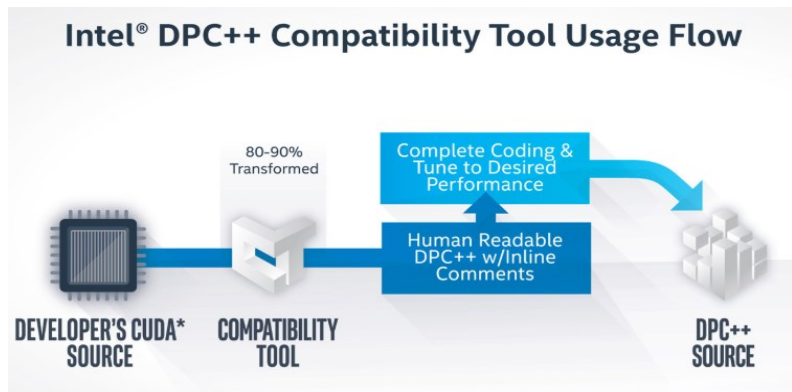


Software beamforming illustration.(Fig source: Lars Grønvold )

Intel's products in software beamforming :

- Core & Gen9 graphics, DG1, Arria 10 & Stratix 10, Intel oneAPI.

# Code Migration

# migration flow



**migration flow**

- The Intel® DPC++ Compatibility Tool assists in migrating your existing CUDA code to Data Parallel C++ (DPC++) code

- DPC++ is based on ISO C++ and incorporates standard SYCL* and community extensions to simplify data parallel programming

- Inline comments help you finish writing and tuning your DPC++ code

| OneAPI Version | Total num of migration place | Success migrated | Need modify | Accuracy |
|---|---|---|---|---|
| beta07 | 84 | 63 | 21 | 75% |
| golden | 84 | 75 | 9 | 89% |

**SUPRA migration summary**

| File Type | *.cpp | *.cu | *.h |
|---|---|---|---|
| File num | 1 | 4 | 23 |

**num. of migrated file**

**Migration Command:** dpct --in-root=./ --out-root=./oneapi  --extra-arg=-Isrc/SupraLib --extra-arg=-Isrc/SupraLib/Beamformer --extra-arg=-Isrc/SupraLib/utilities  --extra-arg=-std=c++11 --extra-arg=-Wno-c++11-narrowing  --extra-arg=-DHAVE_CUDA
./src/SupraLib/Beamformer/ScanConverter.cu ./src/SupraLib/Beamformer/HilbertFirEnvelope.cu
./src/SupraLib/Beamformer/LogCompressor.cu ./src/SupraLib/Beamformer/RxBeamformerCuda.cu ./src/SupraLib/ContainerFactory.cpp

# Migrated code APIs

| Category | oneAPI APIs |
|---|---|
| Memory Management | sycl::malloc_device()<br>sycl::malloc_shared()<br>sycl::free()<br>sycl::malloc_host()<br>sycl::queue.memcpy()<br>Sycl::queue.memset() |
| sycl::queue | dpct::get_current_device().create_queue()<br>dpct::get_default_queue()<br>sycl::queue()<br>sycl::queue.submit()<br>sycl::queue.wait() |
| Math | sycl::sqrt(); sycl::floor(); sycl::fabs(); sycl::round()<br>sycl::max(); sycl::min(); sycl::log10(); sycl::pow() |
| Express Parallel | sycl::nd_item<>; Sycl::nd_range<>; Sycl::range<><br>sycl::id<>; Sycl::nd_item().get_local_range()<br>sycl::nd_item().get_group()<br>sycl::nd_item().get_local_id() |

intel.

# Manually migration example

cuda

```
∨ #ifdef HAVE_CUDA
        cudaEvent_t m_creationEvent;
  #endif
```

**1.** migrate

```
#ifdef HAVE_CUDA
        sycl::event                                        m_creationEvent;
        std::chrono::time_point<std::chrono::steady_clock> m_creationEvent_ct1;
  #endif
```

modify Remove the migrated sycl::event and std::chrono object

cuda

```
cudaSafeCall(cudaStreamCreateWithFlags(&(sm_streams[k]), cudaStreamNonBlocking));
```

**2.** migrate

modify

```
//cudaSafeCall(((sm_streams[k]) = dpct::get_current_device().create_queue()));

sm_streams[k] = new sycl::queue(dpct::get_default_queue_wait().get_context(), dpct::get_default_queue_wait().get_device(), property_list);
```

# Manually migration example



DPCT tool can't migrate CUDA thrust library related code, so it must be rewritten using oneAPI model.

# Migration success example

```
cudaSafeCall(cudaMalloc((void**)&buffer, numBytes));
```

```
cudaSafeCall((buffer = ( uint8_t* )sycl::malloc_device(numBytes, dpct::get_current_device(), dpct::get_default_context()), 0));
```

```
cudaSafeCall(cudaMallocManaged((void**)&buffer, numBytes));
```

```
cudaSafeCall((buffer = ( uint8_t* )sycl::malloc_shared(numBytes, dpct::get_current_device(), dpct::get_default_context()), 0));
```
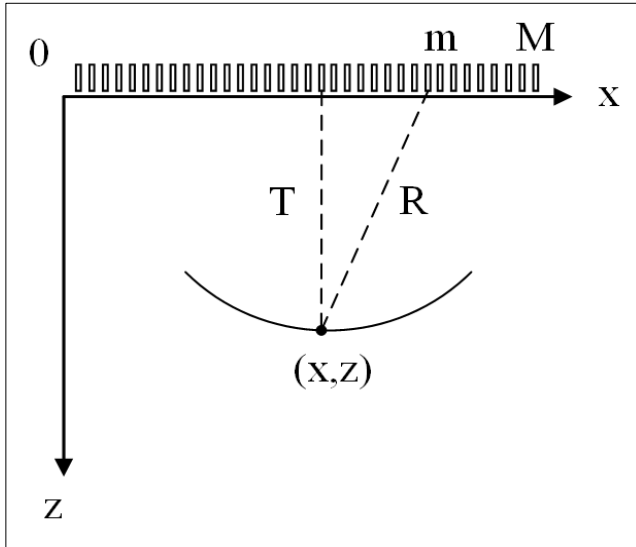
```
cudaSafeCall(cudaMallocHost((void**)&buffer, numBytes));
```

```
cudaSafeCall((buffer = ( uint8_t* )sycl::malloc_host(numBytes, dpct::get_default_context()), 0));
```

Memory allocate related function were successfully migrated.

Table 3  migrated CUDA API summary

intel.

# Beamforming Optimization on GPU

# Beamforming(Delay and Sum) introduction



Geometrical illustration of the pulse-echo process



Delay and sum algorithm illustration

$$\Delta t \; = (T + R)/c_0$$

$c_0$: the speed of ultrasound travel in the body.
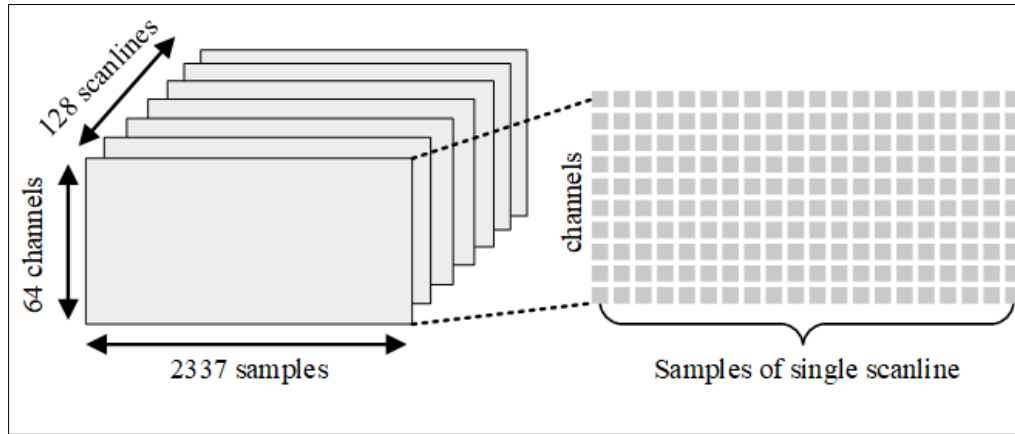
# Beamforming introduction

Suppose pre-beamformed data: 128 scanlines. 64 channels. 2337 samples. It arranged in a 3-D data structure: rf_data[scanline][channel][sample].

Memory required to store pre-beamformed data per image frame is:
(128 * 64 * 2334 * 2) bytes = 36.5 MB
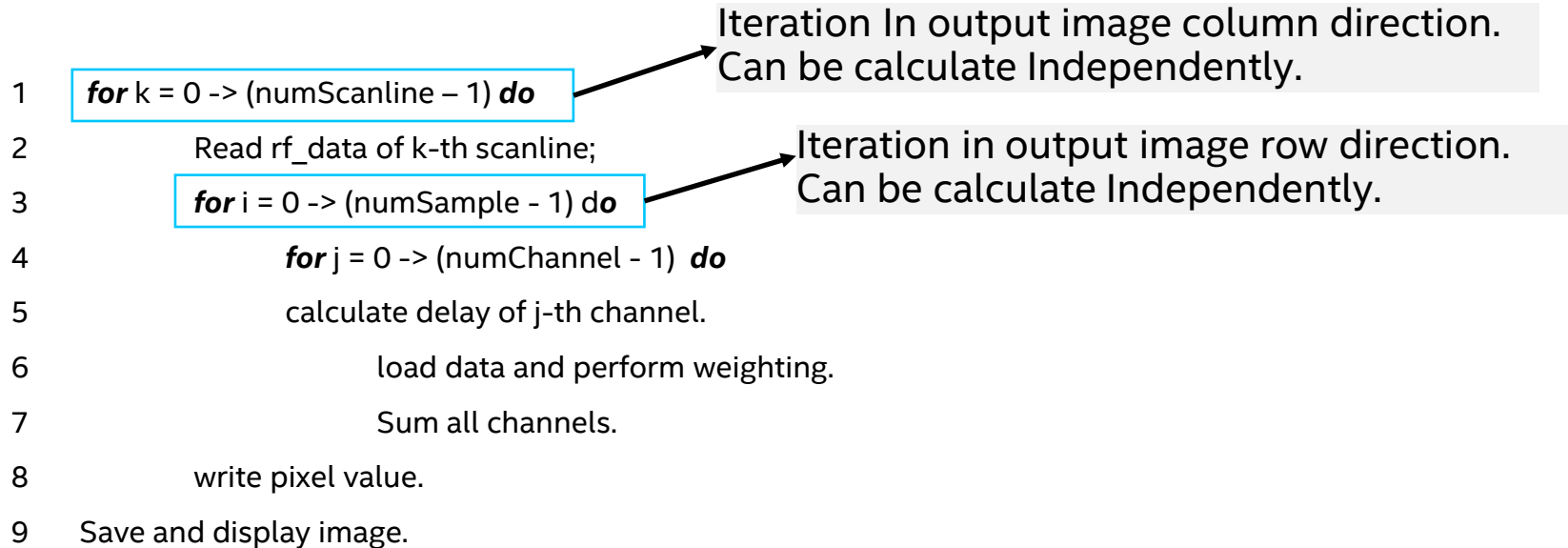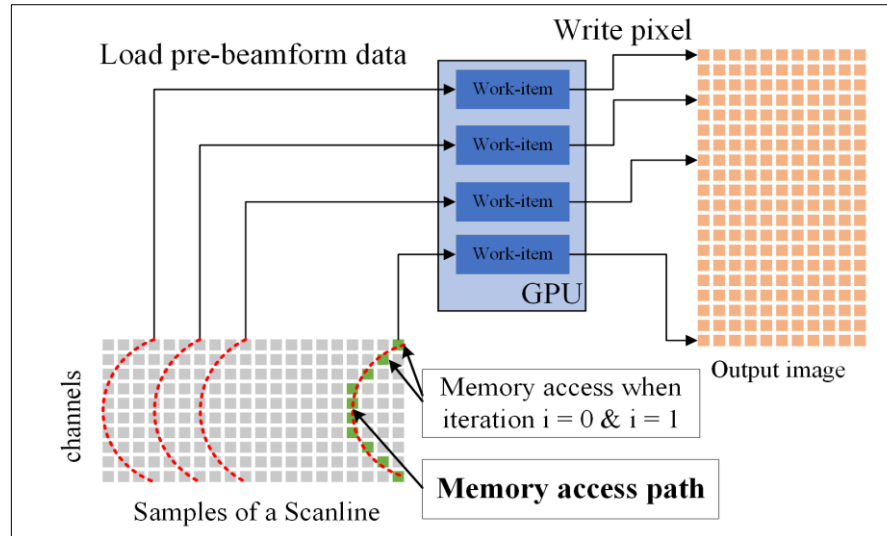


Pre-beamforming data store pattern

# Beamforming algorithm implementation in single thread

Input: 3 dimensional pre-beamformed rf_data[numScanline][numSample][numChannel].

Output: A single frame image.

1       ***for*** k = 0 -> (numScanline – 1) ***do***

2              Read rf_data of k-th scanline;

3             ***for*** i = 0 -> (numSample - 1) d***o***

4                  ***for*** j = 0 -> (numChannel - 1)  ***do***

5                  calculate delay of j-th channel.

6                      load data and perform weighting.

7                    Sum all channels.

8             write pixel value.

9    Save and display image.

Iteration In output image column direction. Can be calculate Independently.

Iteration in output image row direction. Can be calculate Independently.

intel.

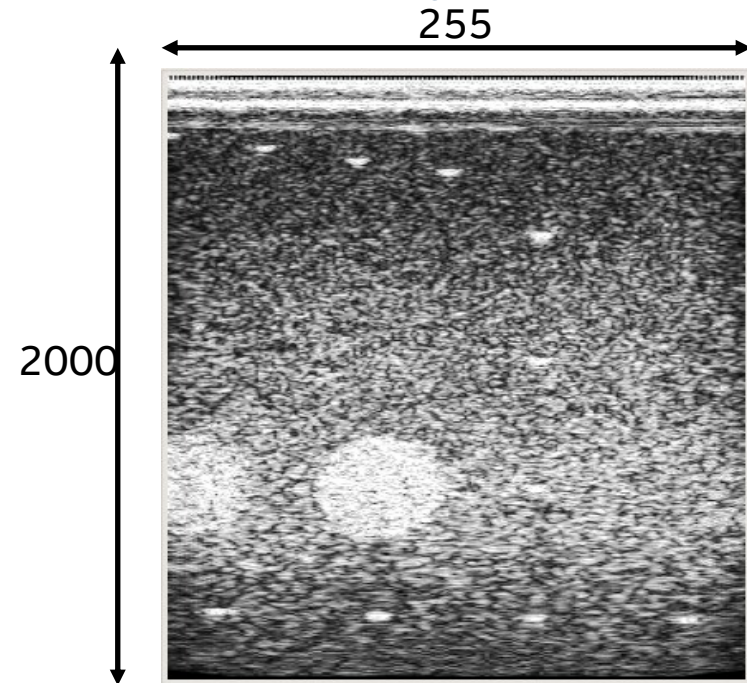# Beamforming implementation in parallel on GPU



Original Beamforming implementation on GPU

# Optimization #1

The optimization is in RxBeamformerCuda.dp.cpp and RxSampleBeamformerDelayAndSum.h file.

Function  rxBeamformingDTSPACEKernel  and sampleBeamfor2D are optimized.



Optimization idea:

**CUDA** :  Each thread calculates a point; every point iterates 64 times.

**oneAPI:** each thread load 2 points in vertical direction, Iterates 8 times.

# Optimization #1

oneAPI code

In RxBeamformerCuda.cu the function been called; the return value is a **float**.

```
#pragma unroll
for (int i = 0; i < row_size; i++) {
    LocationType invMaxElementDistance = 1 / sycl::min(aDT[i], maxElementDistance);
    sInterp[i] = SampleBeamformer::template vec_sampleBeamform2D<interpolateRFlines, RFType, float, LocationType>(txParams, RF, numTransducerElements,
    numReceivedChannels, numTimesteps, x_elemsDT, scanline_x, dirX, dirY, dirZ, aDT[i], d[i], invMaxElementDistance , speedOfSound, dt, additionalOffset,
    windowFunction, mdataGpu);
}
```

In sampleBeamform2D function, calculate single point each call. The for loop at least iterates 64 times.

```
template <bool interpolateRFlines, typename RFType, typename ResultType, typename LocationType>
static ResultType vec_sampleBeamform2D( ScanlineRxParameters3D::TransmitParameters txParams, const RFType* RF, uint32_t numTransducerElements, uint32_t numReceivedChannels,
    uint32_t numTimesteps, const LocationType* x_elemsDT, LocationType scanline_x, LocationType dirX, LocationType dirY, LocationType dirZ, LocationType aDT,
    LocationType depth, LocationType invMaxElementDistance, LocationType speedOfSound, LocationType dt, int32_t additionalOffset,
    const WindowFunctionGpu* __restrict__ windowFunction,const float* mdataGpu)
)
{
    const int VEC_SIZE = 8;
    float sampleAcum = 0.0f;
    float weightAcum = 0.0f;
    int numAdds = 0;
    LocationType initialDelay = txParams.initialDelay;
    uint32_t txScanlineIdx = txParams.txScanlineIdx;

    for (int32_t elemIdxX = txParams.firstActiveElementIndex.x; elemIdxX < txParams.lastActiveElementIndex.x; elemIdxX += VEC_SIZE)
    {
        sycl::vec<int, VEC_SIZE> channelIdx;
        sycl::vec<LocationType, VEC_SIZE> x_elem;          sycl::vec<float, 8>

        #pragma unroll
        for (int i = 0; i < VEC_SIZE; i +=2) {
            channelIdx[i] = (elemIdxX + i) % numReceivedChannels;
            channelIdx[i+1] = (elemIdxX + i + 1) % numReceivedChannels;
            x_elem[i] = x_elemsDT[elemIdxX + i];
            x_elem[i + 1] = x_elemsDT[elemIdxX + i + 1];
        }
        sycl::vec<float, VEC_SIZE> sample;
        sycl::vec<int, VEC_SIZE> mask = (sycl::fabs(x_elem - scanline_x) <= aDT);          Use mask
        /*sycl spec1.2.1 mentioned: true return -1, false return 0*/
        mask *= -1;
        numAdds += utils<int, VEC_SIZE>::add_vec(mask);
```

Source code: supra/src/SupraLib/Beamformer/ RxSampleBeamformerDelayAndSum.h

# Optimization #2

Another optimization in BeamformingNode is directly move into kernel function rather than using function call.

```cpp
sycl::vec<float, VEC_SIZE> weight = windowFunction->get_vec((x_elem - scanline_x) * invMaxElementDistance);

inline sycl::vec<ElementType, VEC_SIZE> get_vec(sycl::vec<float, VEC_SIZE> relativeIndex) const
{
    sycl::vec<float,VEC_SIZE> relativeIndexClamped =
        sycl::min(sycl::max(relativeIndex, -1.0f), 1.0f);
    sycl::vec<float,VEC_SIZE> absoluteIndex =
        m_scale * (relativeIndexClamped + 1.0f);
    sycl::vec<int, VEC_SIZE> int_absoluteIndex = absoluteIndex.convert<int, sycl::rounding_mode::automatic>();

    sycl::vec<float, VEC_SIZE> v(0);
    #pragma unroll
    for(int i = 0 ; i < VEC_SIZE; i++) {
        int index = int_absoluteIndex[i];
        v[i] = m_data[index];       m_data is a private member in class WindowFunctionGpu.
    }
    return v;
}
```

Source code: supra/src/SupraLib/Beamformer/WindowFunction.cpp

Before optimization, fetch data from windowFcuntion->m_data.

intel.

# Optimization #2

Another optimization in BeamformingNode is directly move into kernel function rather than using function call.

```cpp
// use gRawData->getStream copy data to GPU
auto mdataGpu = std::make_shared<Container<float>>(ContainerLocation::LocationGpu, gRawData->getStream(), m_windowFunction->m_data);
```

```cpp
sycl::vec<float, VEC_SIZE> relativeIndex = (x_elem - scanline_x) * invMaxElementDistance;
sycl::vec<float, VEC_SIZE> relativeIndexClamped = sycl::min(sycl::max(relativeIndex, -1.0f), 1.0f);
sycl::vec<float, VEC_SIZE> absoluteIndex = windowFunction->m_scale * (relativeIndexClamped + 1.0f);
sycl::vec<int, VEC_SIZE> absoluteIndex_int = absoluteIndex.convert<int, sycl::rounding_mode::automatic>();
sycl::vec<float, VEC_SIZE> weight;
#pragma unroll
for (int i = 0; i < VEC_SIZE; i += 2 ) {
    weight[i] = mdataGpu[absoluteIndex_int[i]];
    weight[i + 1] = mdataGpu[absoluteIndex_int[i + 1]];
}

//weightAcum += weight;
//numAdds++;
weight *= mask.convert<float, sycl::rounding_mode::automatic>();
weightAcum += utils<float, VEC_SIZE>::add_vec(weight);
```

m_data was copied to mdataGpu befor the queue->submit() call, then mdataGpu was directly passed to kernel function. For data copy, change m_data in WindowFunctionGpu to public.

Source code: supra/src/SupraLib/Beamformer/RxSampleBeamformerDelayAndSum.h

After optimization, fetch data from mdataGpu, mdataGpu was directly pass to kernel function

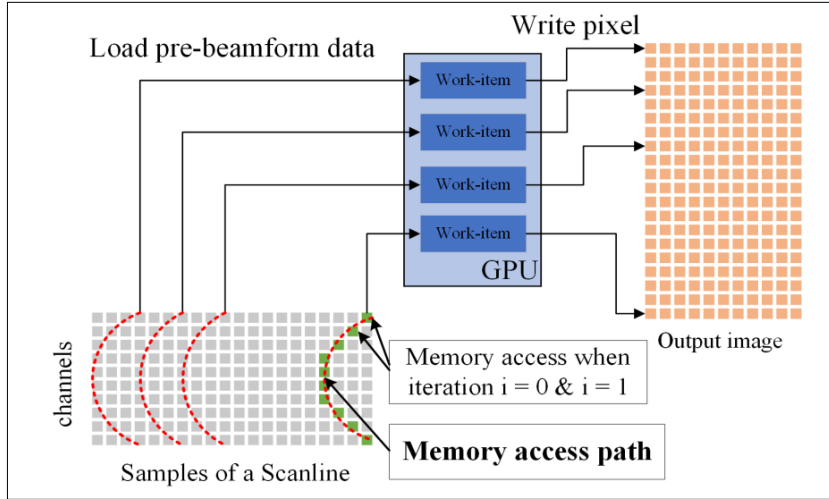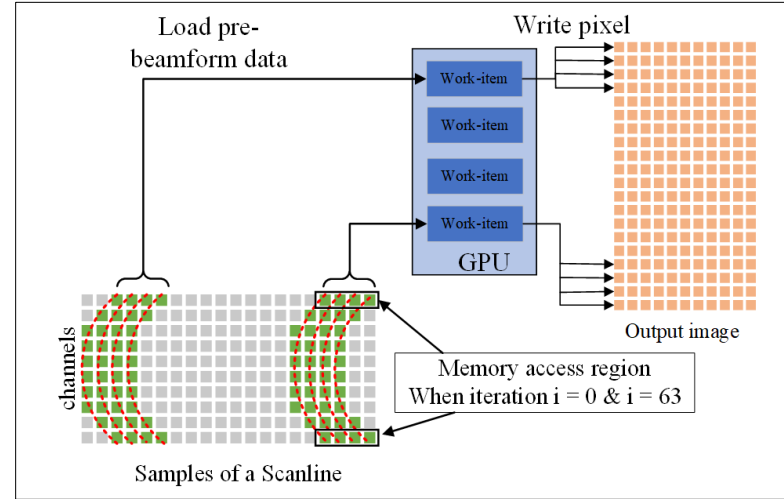intel

# Optimization #2



Before optimization

After optimization

# Using ESIMD to optimize beamforming



Original Beamforming implementation on GPU
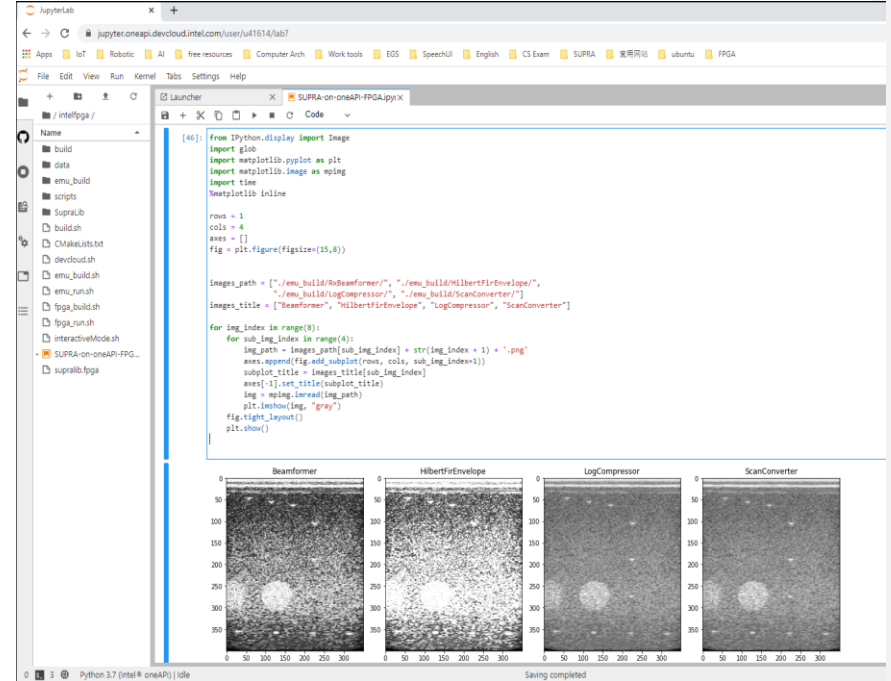
Optimized Beamforming implementation on GPU

# Beamforming implementation on FPGA

# Supra on Intel FPGA Arria 10

| SUPRA Node | oneAPI (ms) UHD630 | oneAPI (ms) Arria 10 |
|---|---|---|
| RxBeamforming | 9.24 | 5.94↓ (Max) |
| HilbertFirEnvelope | 1.50 | 2.61 |
| LogCompressor | 0.27 | 0.34 |
| ScanConverter | 2.65 | 5.66 |
| Total | 13.66 | |



SUPRA on FPGA has been tested on DevCloud, a Jupyter notebook provided to quick build and run.

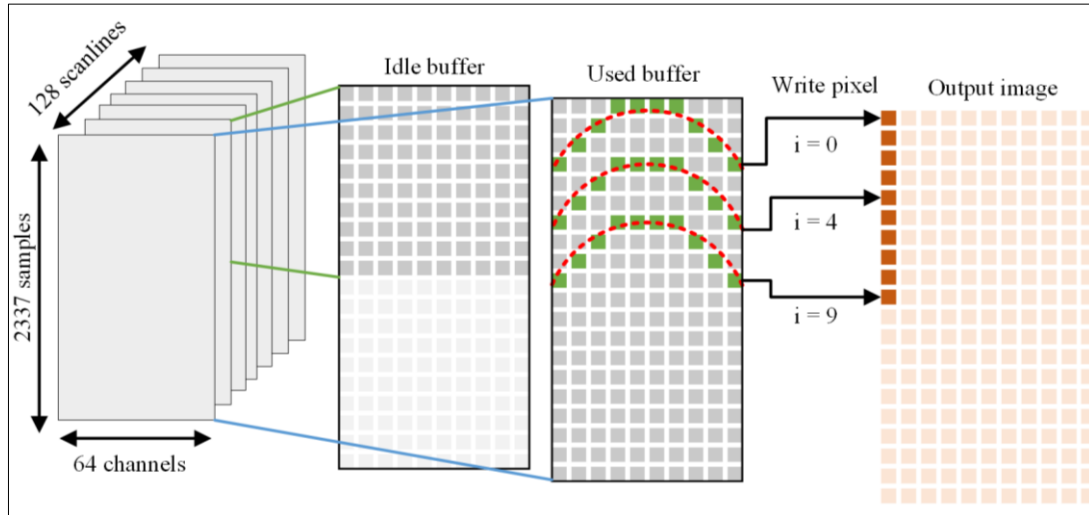Link: https://gitlab.devtools.intel.com/qwang12/ultrasound-emu/-/tree/intelfpga_beta10

# Beamforming on the FPGA



Original pre-beamformed data store pattern

Shuffled pre-beamformed data store pattern

# Beamforming on the FPGA



Beamforming implementation on FPGA



Beamforming algorithm on FPGA

FPGA code: intelfpga-devcloud-golden/SupraLib/Beamformer/ RxBeamformerCuda.dp.cpp

intel.

# Code Sample for FPGA

```
379  template <typename InputType, typename OutputType, typename WeightType, typename IndexType>
380  void scanConvert2D(
381      uint32_t numScanlines,
382      uint32_t numSamples,
383      uint32_t width,
384      uint32_t height,
385      const uint8_t *__restrict__ mask,
386      const IndexType *__restrict__ sampleIdx,
387      const WeightType *__restrict__ weightX,
388      const WeightType *__restrict__ weightY,
389      const InputType *__restrict__ scanlines,
390      OutputType *__restrict__ image,
391      sycl::nd_item<3> item_ct1)
392  {
393      vec2T<uint32_t> pixelPos{
394          item_ct1.get_local_range().get(2) * item_ct1.get_group(2) +
395              item_ct1.get_local_id(2),
396          item_ct1.get_local_range().get(1) * item_ct1.get_group(1) +
397              item_ct1.get_local_id(1)}; //@suppress("Symbol is not resolved")
398                                          //@suppress("Field cannot be resolved")
399
400      if (pixelPos.x < width && pixelPos.y < height)
401      {
402          IndexType pixelIdx = pixelPos.x + pixelPos.y * width;
403          float val = 0;
404          if (mask[pixelIdx])
405          {
406              IndexType sIdx = sampleIdx[pixelIdx];
407              WeightType wX = weightX[pixelIdx];
408              WeightType wY = weightY[pixelIdx];
409
410              val = (1 - wY) * ((1 - wX) * scanlines[sIdx] +
411                      wX * scanlines[sIdx + 1]) +
412                  wY * ((1 - wX) * scanlines[sIdx + numScanlines] +
413                      wX * scanlines[sIdx + 1 + numScanlines]);
414          }
415          image[pixelIdx] = clampCast<OutputType>(val);
416      }
417  }
```

```
424  template <typename InputType, typename OutputType>
425  void buf_fpga_scanConvert2D(
426      uint32_t numScanlines,
427      uint32_t numSamples,
428      uint32_t width,
429      uint32_t height,
430      InputType* scanlines,
431      OutputType* image,
432      sycl::accessor<fpga_data_load, 1, sycl::access::mode::read> fpga_data_load_acc)
433  {
434      uint32_t sIdx;
435      float wX;
436      float wY;
437      float val;
438      int32_t temp;
439      uint32_t buffer_index;
440      uint32_t buffer_1;
441      uint32_t buffer_2;
442      uint32_t buffer_3;
443      int Index_e = 0;

445      // FPGA specific attributes
446      [[intelfpga::max_replicates(4), intelfpga::doublepump]] float buf1[BUFFER_SIZE];
447      [[intelfpga::max_replicates(4), intelfpga::doublepump]] float buf2[BUFFER_SIZE];
448      [[intelfpga::max_replicates(4), intelfpga::doublepump]] float buf3[BUFFER_SIZE];

450      int buf_index = -1;
451      int Index = 0;

453      // preload buffer from global memory.
454      for (uint32_t i = 0; i < BUFFER_SIZE; i++)
455      {
456          buf1[i] = scanlines[i];
457          buf2[i] = scanlines[i + numScanlines];
458      }
```

- oneAPI provides high level language(DPC++) to programming FPGA, which is more flexible, easy to learn, easy to develop, easy to debug.

- To use oneAPI programing for FPGA, Professional knowledge of FPGA is required.

# Results and Performance

# SUPRA GUI and DevCloud usage

intel.

# Intel DevCloud usage – FPGA



Intel DevCloud: https://devcloud.intel.com/oneapi/

intel.

# SUPRA performance on Intel hardware

| SUPRA Node | oneAPI (ms) – UHD630 | Tiger lake Iris Xe | oneAPI - DG1(WA) | oneAPI – Arria 10 |
|---|---|---|---|---|
| RxBeamforming | 9.24 | 4.36 | 3.81 | 5.94 |
| HilbertFirEnvelope | 1.50 | 0.73 | 0.65 | 2.61 |
| LogCompressor | 0.27 | 0.1 | 0.08 | 0.34 |
| ScanConverter | 2.65 | 2.22 | 1.14 | 5.66 |
| Total | 13.66 | 7.41 | 5.68 | 5.94(max) |

For the source code, please refer to:  https://github.com/intel/supra-on-oneapi

For other vendors hardware performance, please refer to:  https://doi.org/10.1007/s11548-018-1750-6

intel.

# Summary

- Unified programing framework/language to implement medical algorithm accelerations on Intel HW
- Samples to implement Ultrasound beamforming on Intel xGPU
- Samples to implement Ultrasound beamforming on Intel FPGA
- Possibility to integrate algorithm acceleration and AI inference on a heterogenous compute system(Intel oneAPI and OpenVINO)
- Future Intel acceleration hardware (xPU) support

intel.

# Thanks for your time!

intel.