# VexCL — a Vector Expression Template Library for OpenCL

Denis Demidov

Supercomputer Center of Russian Academy of Sciences
Kazan Federal University

November 2013, NSCF

# Modern GPGPU frameworks

## CUDA

- Proprietary architecture by NVIDIA
- Requires NVIDIA hardware
- More mature, many libraries

## OpenCL

- Open standard
- Supports wide range of hardware
- Code is much more verbose

# Modern GPGPU frameworks

## CUDA

- Proprietary architecture by NVIDIA
- Requires NVIDIA hardware
- More mature, many libraries

- *Kernels are compiled to PTX together with host program*

## OpenCL

- Open standard
- Supports wide range of hardware
- Code is much more verbose

- *Kernels are compiled at runtime, adding an initialization overhead*

- The latter distinction is usually considered to be an OpenCL drawback.
- But it also allows us to generate more efficient kernels at runtime!
  - VexCL takes care of this part.

# VexCL — a vector expression template library for OpenCL

- https://github.com/ddemidov/vexcl

- Created for ease of C++ based OpenCL development.
  - □ Convenient notation for vector expressions.
  - □ OpenCL JIT code generation.
- The source code is publicly available under MIT license.
- *This is not a* C++ *bindings library!*
  - □ VexCL works on top of Khronos C++ bindings for OpenCL.

# Hello OpenCL: vector sum

## Vector sum

- $A$, $B$, and $C$ are large vectors.
- Compute $C = A + B$.

## Overview of OpenCL solution

1. Initialize OpenCL context on supported device.
2. Allocate memory on the device.
3. Transfer input data to device.
4. Run your computations on the device.
5. Get the results from the device.

### 1. Query platforms

```
1  std :: vector<cl::Platform> platform;
2  cl :: Platform::get(&platform);
3
4  if ( platform.empty() )
5      throw std::runtime_error("OpenCL platforms not found.");
```

### 2. Get first available GPU device

```cpp
6   cl :: Context context;
7   std :: vector<cl::Device> device;
8   for(auto p = platform.begin(); device.empty() && p != platform.end(); p++) {
9       std :: vector<cl::Device> dev;
10      try {
11          p−>getDevices(CL_DEVICE_TYPE_GPU, &dev);
12          for(auto d = dev.begin(); device.empty() && d != dev.end(); d++) {
13              if (!d−>getInfo<CL_DEVICE_AVAILABLE>()) continue;
14              device. push_back(∗d);
15              context = cl :: Context(device);
16          }
17      } catch (...) {
18          device. clear ();
19      }
20  }
21  if (device.empty()) throw std::runtime_error("GPUs not found");
```

### 3. Create kernel source

```
22  const char source[] =
23      "kernel void add(\n"
24      "        uint n,\n"
25      "        global const float *a,\n"
26      "        global const float *b,\n"
27      "        global float *c\n"
28      "        )\n"
29      "{\n"
30      "    uint i = get_global_id(0);\n"
31      "    if (i < n) {\n"
32      "        c[i] = a[i] + b[i];\n"
33      "    }\n"
34      "}\n";
```

# Hello OpenCL: vector sum

### 4. Compile kernel

```
35   cl :: Program program(context, cl::Program::Sources(
36              1, std :: make_pair(source, strlen (source))
37              ));
38   try {
39       program.build(device);
40   } catch (const cl::Error&) {
41       std :: cerr
42           << "OpenCL compilation error" << std::endl
43           << program.getBuildInfo<CL_PROGRAM_BUILD_LOG>(device[0])
44           << std::endl;
45       return 1;
46   }
47   cl :: Kernel add_kernel = cl :: Kernel(program, "add");
```

### 5. Create command queue

```
48   cl :: CommandQueue queue(context, device[0]);
```

### 6. Prepare input data, transfer it to device

```
49   const size_t N = 1 << 20;
50   std :: vector<float> a(N, 1), b(N, 2), c(N);
51
52   cl :: Buffer  A(context, CL_MEM_READ_ONLY | CL_MEM_COPY_HOST_PTR,
53           a. size () * sizeof(float ),  a.data ());
54
55   cl :: Buffer  B(context, CL_MEM_READ_ONLY | CL_MEM_COPY_HOST_PTR,
56           b. size () * sizeof(float ),  b.data ());
57
58   cl :: Buffer  C(context, CL_MEM_READ_WRITE,
59           c. size () * sizeof(float ));
```

### 7. Set kernel arguments

```
60   add_kernel.setArg(0, N);
61   add_kernel.setArg(1, A);
62   add_kernel.setArg(2, B);
63   add_kernel.setArg(3, C);
```

### 8. Launch kernel

```
64   queue.enqueueNDRangeKernel(add_kernel, cl::NullRange, N, cl::NullRange);
```

### 9. Get result back to host

```
65   queue.enqueueReadBuffer(C, CL_TRUE, 0, c.size() * sizeof(float), c.data());
66   std::cout << c[42] << std::endl; // Should get '3' here.
```

## Hello VexCL: vector sum

**Much shorter!**

```
1   std :: cout << 3 << std::endl;
```

# Hello VexCL: vector sum

### Get all available GPUs:

```
1  vex::Context ctx( vex::Filter::Type(CL_DEVICE_TYPE_GPU) );
2  if ( !ctx ) throw std::runtime_error("GPUs not found");
```

### Prepare input data, transfer it to device:

```
3  std::vector<float> a(N, 1), b(N, 2), c(N);
4  vex::vector<float> A(ctx, a);
5  vex::vector<float> B(ctx, b);
6  vex::vector<float> C(ctx, N);
```

### Launch kernel, get result back to host:

```
7  C = A + B;
8  vex::copy(C, c);
9  std::cout << c[42] << std::endl;
```
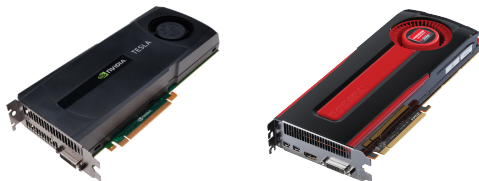
## Initialization

- Multi-device and multi-platform computations are supported.
- VexCL context is initialized from combination of device filters.
- Device filter is a boolean functor acting on **const** cl::Device&.

### Initialize VexCL context on selected devices

```
1   vex::Context ctx( vex::Filter :: All );
```

## Initialization

- Multi-device and multi-platform computations are supported.
- VexCL context is initialized from combination of device filters.
- Device filter is a boolean functor acting on **const** cl::Device&.

### Initialize VexCL context on selected devices

```
1   vex::Context ctx( vex::Filter::Type(CL_DEVICE_TYPE_GPU) );
```

## Initialization

- Multi-device and multi-platform computations are supported.
- VexCL context is initialized from combination of device filters.
- Device filter is a boolean functor acting on **const** cl::Device&.

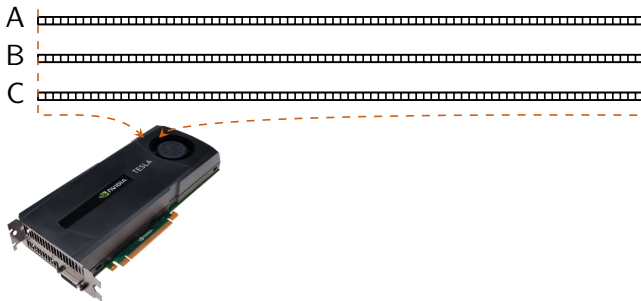### Initialize VexCL context on selected devices

```
1  vex::Context ctx(
2      vex::Filter::Type(CL_DEVICE_TYPE_GPU) &&
3      vex::Filter::Platform("AMD")
4      );
```

## Initialization

- Multi-device and multi-platform computations are supported.
- VexCL context is initialized from combination of device filters.
- Device filter is a boolean functor acting on **const** cl::Device&.

### Initialize VexCL context on selected devices

```
1  vex::Context ctx(
2      vex::Filter::Type(CL_DEVICE_TYPE_GPU) &&
3      [](const cl::Device &d) {
4          return d.getInfo<CL_DEVICE_GLOBAL_MEM_SIZE>() >= 4_GB;
5      });
```

## Memory and work splitting
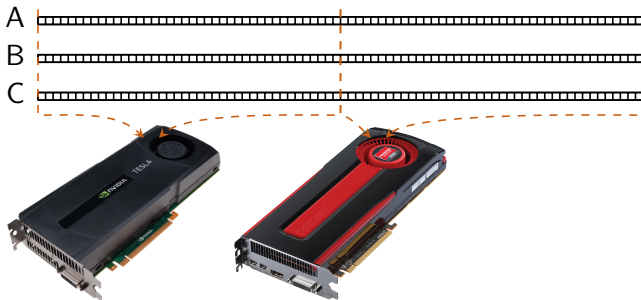
### Hello VexCL example

```
1  vex::Context ctx( vex::Filter::Name("Tesla") );
2
3  vex::vector<float> A(ctx, N); A = 1;
4  vex::vector<float> B(ctx, N); B = 2;
5  vex::vector<float> C(ctx, N);
6
7  C = A + B;
```

# Memory and work splitting
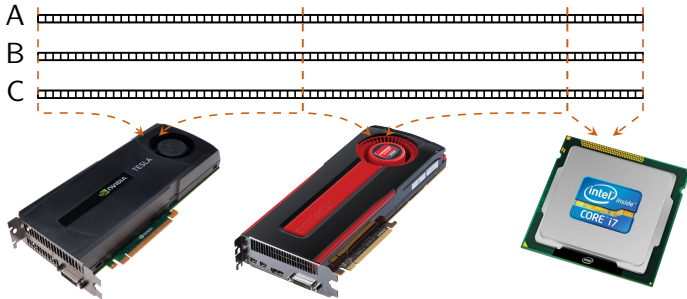
## Hello VexCL example

```
1  vex::Context ctx( vex::Filter::Type(CL_DEVICE_TYPE_GPU) );
2
3  vex::vector<float> A(ctx, N); A = 1;
4  vex::vector<float> B(ctx, N); B = 2;
5  vex::vector<float> C(ctx, N);
6
7  C = A + B;
```

### Hello VexCL example

```
1   vex::Context ctx( vex::Filter::DoublePrecision );
2
3   vex::vector<float> A(ctx, N); A = 1;
4   vex::vector<float> B(ctx, N); B = 2;
5   vex::vector<float> C(ctx, N);
6
7   C = A + B;
```

## Copies between host and device memory

```
1  vex::vector<double> X;
2  std::vector<double> x;
3  double c_array[100];
```

### Simple copies

```
1  vex::copy(X, x); // From device to host.
2  vex::copy(x, X); // From host to device.
```

### STL-like range copies

```
1  vex::copy(X.begin(), X.end(), x.begin());
2  vex::copy(X.begin(), X.begin() + 100, x.begin());
3  vex::copy(c_array, c_array + 100, X.begin());
```

### Inspect or set single element (*slow*)

```
1  assert(x[42] == X[42]);
2  X[0] = 0;
```

## What vector expressions are supported?

- All vectors in an expression have to be *compatible*:
  - □ Have same size
  - □ Located on same devices

- What may be used:

  - □ Vectors and scalars
  - □ Arithmetic, logical operators
  - □ Built-in OpenCL functions
  - □ User-defined functions
  - □ Random number generators

  - □ Slicing and permutations
  - □ Reduction (sum, min, max)
  - □ Stencil operations
  - □ Sparse matrix – vector products
  - □ Fast Fourier Transform

```
1  vex::vector<double> x(ctx, n), y(ctx, n);
2
3  x = (2 * M_PI / n) * vex::element_index();
4  y = pow(sin(x), 2.0) + pow(cos(x), 2.0);
```
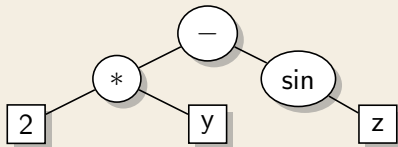
## Builtin operations and functions

This expression:

```
1   x = 2 * y − sin(z);
```

- Define VEXCL_SHOW_KERNELS to see the generated code.

... results in this kernel:

```
1   kernel void vexcl_vector_kernel(
2       ulong n,
3       global double ∗ res,
4       int prm1,
5       global double ∗ prm2,
6       global double ∗ prm3
7   )
8   {
9       for(size_t idx = get_global_id(0); idx < n; idx += get_global_size(0)) {
10          res[idx] = ( ( prm1 ∗ prm2[idx] ) − sin( prm3[idx] ) );
11      }
12  }
```

# Element indices

- vex :: element_index(**size_t** offset = 0) returns index of an element inside a vector.
  - □ The numbering starts with offset and is continuous across devices.

### Linear function:

```
1  vex :: vector<double> X(ctx, N);
2  double x0 = 0, dx = 1e−3;
3  X = x0 + dx * vex::element_index();
```

### Single period of sine function:

```
1  X = sin(2 * M_PI * vex::element_index() / N);
```

# User-defined functions

- Users may define functions to be used in vector expressions.
- There are two options for doing this:
  - □ Provide function body in a string.
  - □ Provide generic C++ functor.
- Once defined, user functions are used in the same way as builtin functions.

# 1. Provide function body in a string

- Choose function name
- Specify function signature
- Provide function body

Defining a function:

```
1  VEX_FUNCTION( sqr, double(double, double), "return prm1 * prm1 + prm2 * prm2;" );
```

Using a function:

```
1  Z = sqrt( sqr(X, Y) );
```

## 2. Provide generic functor

Functor definition:

```
1  struct sqr_functor {
2      template <class T>
3      T operator(const T &x, const T &y) const {
4          return x * x + y * y;
5      }
6  };
```

# 2. Provide generic functor

Functor definition:

```
1  struct sqr_functor {
2      template <class T>
3      T operator(const T &x, const T &y) const {
4          return x * x + y * y;
5      }
6  };
```

Generate VexCL function:

```
1  using vex::generator::make_function;
2  auto sqr = make_function<double(double,double)>( sqr_functor() );
```

# 2. Provide generic functor

### Functor definition:

```
1  struct sqr_functor {
2      template <class T>
3      T operator(const T &x, const T &y) const {
4          return x * x + y * y;
5      }
6  };
```

### Generate VexCL function:

```
1  using vex::generator::make_function;
2  auto sqr = make_function<double(double,double)>( sqr_functor() );
```

### Boost.Phoenix lambdas *are* generic functors:

```
1  using namespace boost::phoenix::arg_names;
2  auto sqr = make_function<double(double,double)>( arg1 * arg1 + arg2 * arg2 );
```
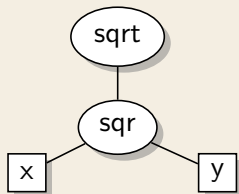
## User functions are translated to OpenCL functions

```
1   Z = sqrt( sqr(X, Y) );
```

### . . . gets translated to:

```
1   double func1(double prm1, double prm2) {
2       return prm1 * prm1 + prm2 * prm2;
3   }
4
5   kernel void vexcl_vector_kernel(
6       ulong n,
7       global double * res,
8       global double * prm1,
9       global double * prm2
10  )
11  {
12      for(size_t idx = get_global_id(0); idx < n; idx += get_global_size(0)) {
13          res[idx] = sqrt( func1( prm1[idx], prm2[idx] ) );
14      }
15  }
```
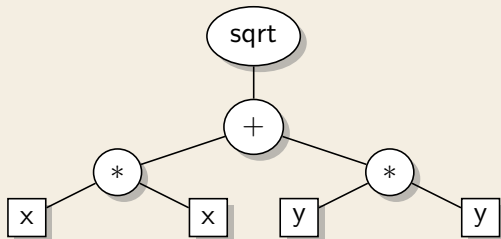
## Functions may be not only convenient, but also effective

### Same example without using a function:

```
1   Z = sqrt( X * X + Y * Y );
```

### . . . gets translated to:

```
1   kernel void vexcl_vector_kernel(
2       ulong n,
3       global double * res,
4       global double * prm1,
5       global double * prm2,
6       global double * prm3,
7       global double * prm4
8   )
9   {
10      for(size_t idx = get_global_id(0); idx < n; idx += get_global_size(0)) {
11          res[idx] = sqrt( ( ( prm1[idx] * prm2[idx] ) + ( prm3[idx] * prm4[idx] ) ) );
12      }
13  }
```

## Tagged terminals

- Programmer may help VexCL to recognize same terminals by tagging them:

### Like this:

```
1  using vex::tag;
2  Z = sqrt(tag<1>(X) * tag<1>(X) +
3          tag<2>(Y) * tag<2>(Y));
```

### or, equivalently:

```
1  auto x = tag<1>(X);
2  auto y = tag<2>(Y);
3  Z = sqrt(x * x + y * y);
```
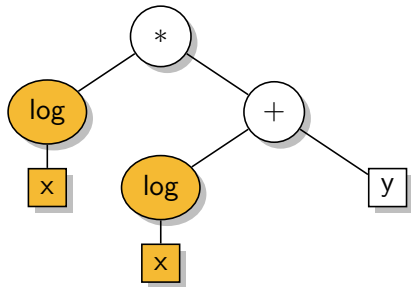
```
1  kernel void vexcl_vector_kernel(
2      ulong n,
3      global double * res,
4      global double * prm1,
5      global double * prm2
6  )
7  {
8      for(size_t idx = get_global_id(0); idx < n; idx += get_global_size(0)) {
9          res[idx] = sqrt( ( ( prm1[idx] * prm1[idx] ) + ( prm2[idx] * prm2[idx] ) ) );
10     }
11 }
```

## Reusing intermediate results

- Some expressions may have several inclusions of the same subexpression:

1  $Z = \log(X) * (\log(X) + Y);$

  - $\log(X)$ will be computed twice here.
  - One could tag $X$ and hope that OpenCL compiler is smart enough...

## Temporaries

- But it is also possible to introduce a temporary variable explicitly:

```
1  auto tmp = vex::make_temp<1>( log(X) );
2  Z = tmp * (tmp + Y);
```

```
1  kernel void vexcl_vector_kernel(
2      ulong n,
3      global double * res,
4      global double * prm1,
5      global double * prm2
6  )
7  {
8      for(size_t idx = get_global_id(0); idx < n; idx += get_global_size(0)) {
9          double temp1 = log( prm1[idx] );
10         res[idx] = ( temp1 * ( temp1 + prm2[idx] ) );
11     }
12 }
```

## Slicing    (*Single-device contexts*)

- When working with dense multidimensional matrices, it is general practice to store those in continuous arrays.
    - □ An instance of vex::slicer <NDIM> class allows to access sub-blocks of such matrix.

n-by-n matrix and a slicer:

```
1  vex::vector<double> x(ctx, n * n);
2  vex::slicer <2> slice(vex::extents[n][n]); // Can be used with any vector of appropriate size
```

## Slicing      (*Single-device contexts*)

- When working with dense multidimensional matrices, it is general practice to store those in continuous arrays.
  - An instance of vex::slicer<NDIM> class allows to access sub-blocks of such matrix.

n-by-n matrix and a slicer:

```
1  vex::vector<double> x(ctx, n * n);
2  vex::slicer<2> slice(vex::extents[n][n]);  // Can be used with any vector of appropriate size
```

Access row or column of the matrix:

```
3  using vex::_;
4  y = slice[42](x);        // 42nd row
5  y = slice[_][42](x);     // 42nd column
6  slice[_][10](x) = y;     // Slices are writable
```

## Slicing (*Single-device contexts*)

- When working with dense multidimensional matrices, it is general practice to store those in continuous arrays.
  - An instance of vex:: slicer <NDIM> class allows to access sub-blocks of such matrix.

### n-by-n matrix and a slicer:

```
1   vex::vector<double> x(ctx, n * n);
2   vex:: slicer <2> slice(vex::extents[n][n]); // Can be used with any vector of appropriate size
```

### Access row or column of the matrix:

```
3   using vex::_;
4   y = slice [42](x);        // 42nd row
5   y = slice [_][42]( x);    // 42nd column
6   slice [_][10]( x) = y;    // Slices are writable
```

### Use ranges to select sub-blocks:

```
7   using vex::range;
8   z = slice [range(0, 2, n )][range(10, 20)](x);
```

- vex::permutation() function takes arbitrary (integral valued) vector expression and returns permutation functor:

## Permutations  (*Single-device contexts*)

- $vex::permutation()$ function takes arbitrary (integral valued) vector expression and returns permutation functor:

Index-based permutation:

```
1  vex::vector<size_t> I(ctx, N);
2  I = N − 1 − vex::element_index();
3  auto reverse = vex::permutation(I);
4  y = reverse(x);
```

## Permutations  (*Single-device contexts*)

- vex::permutation() function takes arbitrary (integral valued) vector expression and returns permutation functor:

### Index-based permutation:

```
1   vex::vector<size_t> I(ctx, N);
2   I = N − 1 − vex::element_index();
3   auto reverse = vex::permutation(I);
4   y = reverse(x);
```

### Expression-based permutation:

```
1   auto reverse = vex::permutation(N − 1 − vex::element_index());
2   y = reverse(x);
```

## Permutations  (*Single-device contexts*)

- vex::permutation() function takes arbitrary (integral valued) vector expression and returns permutation functor:

### Index-based permutation:

```
1  vex::vector<size_t> I(ctx, N);
2  I = N − 1 − vex::element_index();
3  auto reverse = vex::permutation(I);
4  y = reverse(x);
```

### Expression-based permutation:

```
1  auto reverse = vex::permutation(N − 1 − vex::element_index());
2  y = reverse(x);
```

### Permutations are writable:

```
1  reverse(y) = x;
```

## Random number generation

- VexCL provides implementation[1] of *counter-based* random number generators from Random123[2] suite.
  - The generators are *stateless*; mixing functions are applied to element indices.
  - Implemented families: Threefry and Philox.
- vex::Random<T,G> — uniform distribution.
- vex::RandomNormal<T,G> — normal distribution.

### Monte Carlo $\pi$:

```
1  vex::Random<double, vex::random::threefry> rnd;
2
3  x = 2 * rnd(vex::element_index(), std::rand()) − 1;
4  y = 2 * rnd(vex::element_index(), std::rand()) − 1;
5
6  vex::Reductor<size_t, vex::SUM> sum(ctx);
7  double pi = 4.0 * sum( (x * x + y * y) < 1 ) / n;
```

[1]Contributed by Pascal Germroth ⟨pascal@ensieve.org⟩
[2]D E Shaw Research, http://www.deshawresearch.com/resources_random123.html

## Reductions

- Class vex::Reductor<T, kind> allows to reduce arbitrary *vector expression* to a single value of type T.
- Supported reduction kinds: SUM, MIN, MAX

### Inner product

```
1  vex::Reductor<double, vex::SUM> sum(ctx);
2  double s = sum(x * y);
```

### Number of elements in x between 0 and 1

```
1  vex::Reductor<size_t, vex::SUM> sum(ctx);
2  size_t n = sum( (x > 0) && (x < 1) );
```

### Maximum distance from origin

```
1  vex::Reductor<double, vex::MAX> max(ctx);
2  double d = max( sqrt(x * x + y * y) );
```

## Sparse matrix – vector products    *(Additive expressions only)*

- Class vex::SpMat<T> holds representation of a sparse matrix on compute devices.
- Constructor accepts matrix in common CRS format:
  - row indices, columns and values of nonzero entries.

### Construct matrix

```
1   vex::SpMat<double> A(ctx, n, n, row.data(), col.data(), val.data());
```

### Compute residual value

```
2   // vex::vector<double> u, f, r;
3   r = f − A * u;
4   double res = max( fabs(r) );
```

- SpMV may only be used in additive expressions:
  - □ Needs data exchange between compute devices.
  - □ Impossible to implement with single kernel.

- This restriction may be lifted for single-device contexts:

```
r = f − vex::make_inline(A ∗ u);
double res = max( fabs(r) );
```
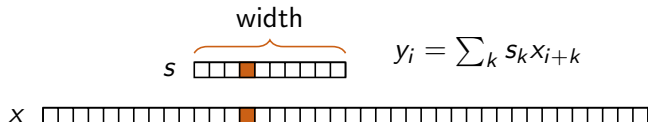
- SpMV may only be used in additive expressions:
  - □ Needs data exchange between compute devices.
  - □ Impossible to implement with single kernel.
- This restriction may be lifted for single-device contexts:

```
r = f − vex::make_inline(A ∗ u);
double res = max( fabs(r) );
```

Do not store intermediate results:

```
double res = max( fabs( f − vex::make_inline(A ∗ u) ) );
```

# Simple stencil convolutions   *(Additive expressions only)*



width

$s$

$$y_i = \sum_k s_k x_{i+k}$$

$x$

- Simple stencil is based on a 1D array, and may be used for:
  - Signal filters (e.g. averaging)
  - Differential operators with constant coefficients
  - . . .

## Moving average with 5-points window

```cpp
std :: vector<double> sdata(5, 0.2);
vex :: stencil <double> s(ctx, sdata, 2 /* center */);

y = x * s;
```

## User-defined stencil operators (*Additive expressions only*)

- Define efficient arbitrary stencil operators:
  - □ Return type
  - □ Stencil dimensions (width and center)
  - □ Function body
  - □ Queue list

### Example: nonlinear operator

$$y_i = x_i + (x_{i-1} + x_{i+1})^3$$

### Implementation

```
1  VEX_STENCIL_OPERATOR(custom_op, double, 3/*width*/, 1/*center*/,
2      "double t = X[-1] + X[1];\n"
3      "return X[0] + t * t * t;",
4      ctx);
5
6  y = custom_op(x);
```

## Fast Fourier Transform   (*Single-device contexts*)

- VexCL provides FFT implementation[3]:
    - Currently only single-device contexts are supported
    - Arbitrary vector expressions as input
    - Multidimensional transforms
    - Arbitrary sizes

### Solve Poisson equation with FFT:

```
1  vex::FFT<double, cl_double2> fft(ctx, n);
2  vex::FFT<cl_double2, double> ifft(ctx, n, vex::inverse);
3
4  vex::vector<double>  rhs(ctx, n), u(ctx, n), K(ctx, n);
5  // ... initialize vectors ...
6
7  u = ifft ( K * fft (rhs) );
```

---

[3]Contributed by Pascal Germroth ⟨pascal@ensieve.org⟩

## Multivectors

- vex::multivector<T,N> holds N instances of equally sized vex::vector<T>
- Supports all operations that are defined for vex::vector<>.
- Transparently dispatches the operations to the underlying components.
- vex::multivector::**operator**()(**size_t** k) returns k-th component.

```
1   vex::multivector<double, 2> X(ctx, N), Y(ctx, N);
2   vex::Reductor<double, vex::SUM> sum(ctx);
3   vex::SpMat<double> A(ctx, ... );
4   std::array<double, 2> v;
5
6   // ...
7
8   X = sin(v * Y + 1);                  // X(k) = sin(v[k] * Y(k) + 1);
9   v = sum( between(0, X, Y) );         // v[k] = sum( between( 0, X(k), Y(k) ) );
10  X = A * Y;                           // X(k) = A * Y(k);
```

## Multiexpressions

- Sometimes an operation cannot be expressed with simple multivector arithmetics.

### Example: rotate 2D vector by an angle

$$y_0 = x_0 \cos \alpha - x_1 \sin \alpha,$$
$$y_1 = x_0 \sin \alpha + x_1 \cos \alpha.$$

- Multiexpression is a tuple of normal vector expressions
- Its assignment to a multivector is functionally equivalent to component-wise assignment, but results in a single kernel launch.

## Multiexpressions

- Multiexpressions may be used with multivectors:

```
1  // double alpha;
2  // vex::multivector<double,2> X, Y;
3
4  Y = std::tie ( X(0) * cos(alpha) − X(1) * sin(alpha),
5                 X(0) * sin(alpha) + X(1) * cos(alpha)  );
```

- and with tied vectors:

```
1  // vex::vector<double> alpha;
2  // vex::vector<double> odlX, oldY, newX, newY;
3
4  vex::tie (newX, newY) = std::tie( oldX * cos(alpha) − oldY * sin(alpha),
5                                    oldX * sin(alpha) + oldY * cos(alpha)  );
```

## A multiexpression results in a single kernel

```
1  auto x0 = tag<0>( X(0) );
2  auto x1 = tag<1>( X(1) );
3  auto ca = tag<2>( cos(alpha) );
4  auto sa = tag<3>( sin(alpha) );
5
6  Y = std::tie(x0 * ca − x1 * sa, x0 * sa + x1 * ca);
```

```
1  kernel void vexcl_multivector_kernel(ulong n,
2      global double * res1, global double * res2,
3      global double * prm1, double prm2,
4      global double * prm3, double prm4
5  )
6  {
7      for(size_t idx = get_global_id(0); idx < n; idx += get_global_size(0)) {
8          double buf1 = ( ( prm1[idx] * prm2 ) − ( prm3[idx] * prm4 ) );
9          double buf2 = ( ( prm1[idx] * prm4 ) + ( prm3[idx] * prm2 ) );
10
11         res1[idx] = buf1;
12         res2[idx] = buf2;
13     }
14 }
```

# Parameter study for the Lorenz attractor system

## Lorenz attractor system

$$\dot{x} = -\sigma (x - y),$$
$$\dot{y} = Rx - y - xz,$$
$$\dot{z} = -bz + xy.$$

- Let's solve large number of Lorenz systems, each for a different value of $R$.
- Let's use VexCL and Boost.odeint for that.



Lorenz attractor trajectory

## Using Boost.odeint

ODE in general:

$$\frac{\mathrm{d}x}{\mathrm{d}t} = \dot{x} = f(x, t), \qquad x(0) = x_0.$$

Using Boost.odeint:

1. Define state type (what is $x$?)
2. Provide system function (define $f$)
3. Choose integration method
4. Integrate over time

## Naive implementation

### 1. State type

```
1  typedef vex::multivector<double, 3> state_type;
```

### 2. System functor

```
2  struct lorenz_system {
3      const vex::vector<double> &R;
4      lorenz_system(const vex::vector<double> &R ) : R(R) { }
5
6      void operator()(const state_type &x, state_type &dxdt, double t) {
7          dxdt = std::tie ( sigma * ( x(1) − x(0) ),
8                            R * x(0) − x(1) − x(0) * x(2),
9                            x(0) * x(1) − b * x(2)   );
10     }
11 };
```

# Naive implementation

### 3. Stepper (4th order Runge-Kutta)

```
12    odeint :: runge_kutta4<
13          state_type  /*state*/,        double  /*value*/,
14          state_type  /*derivative*/,  double  /*time*/,
15          odeint :: vector_space_algebra,  odeint :: default_operations
16          > stepper;
```

### 4. Integration

```
17    vex :: multivector<double,3> X(ctx, n);
18    vex :: vector<double> R(ctx, n);
19
20    X = 10;
21    R = Rmin + vex::element_index() * ((Rmax − Rmin) / (n − 1));
22
23    odeint :: integrate_const (stepper, lorenz_system(R), X, 0.0, t_max, dt);
```

## CUBLAS implementation

- CUBLAS is a highly optimized BLAS implementation from NVIDIA.
- Its disadvantage is that it has fixed number of kernels/functions.
- Hence, linear combinations (used internally by odeint):

$$x_0 = \alpha_1 x_1 + \alpha_2 x_2 + \cdots + \alpha_n x_n$$

are implemented as:

```
cublasDcopy (...);        // x_0 = x_1
cublasDscal (...);        // x_0 = α_1 * x_0
cublasDaxpy (...);        // x_0 = x_0 + α_2 * x_2
...
cublasDaxpy (...);        // x_0 = x_0 + α_n * x_n
```

# Thrust implementation

- It is possible to fuse linear combination kernels with Thrust:
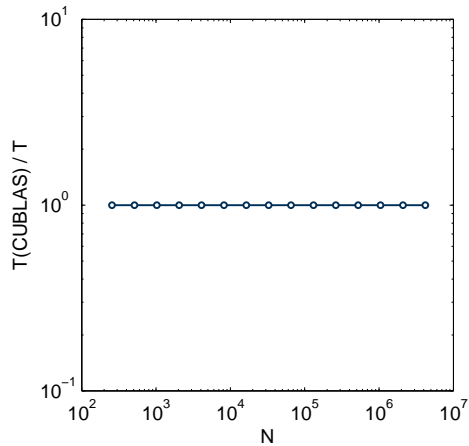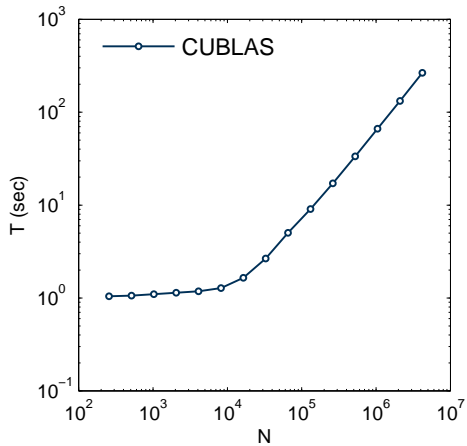
### Thrust

```
1   struct scale_sum2 {
2       const double a1, a2;
3       scale_sum2(double a1, double a2) : a1(a1), a2(a2) { }
4       template<class Tuple>
5       __host__ __device__ void operator()(Tuple t) const {
6           thrust::get<0>(t) = a1 * thrust::get<1>(t) + a2 * thrust::get<2>(t);
7       }
8   };
9
10  thrust::for_each(
11          thrust::make_zip_iterator(
12              thrust::make_tuple( x0.begin(), x1.begin(), x2.begin() )
13              ),
14          thrust::make_zip_iterator(
15              thrust::make_tuple( x0.end(), x1.end(), x2.end() )
16              ),
17          scale_sum2(a1, a2)
18          );
```
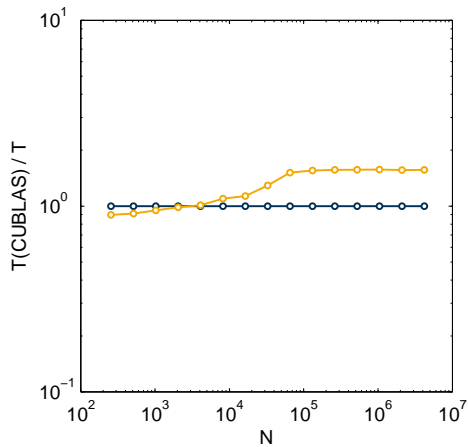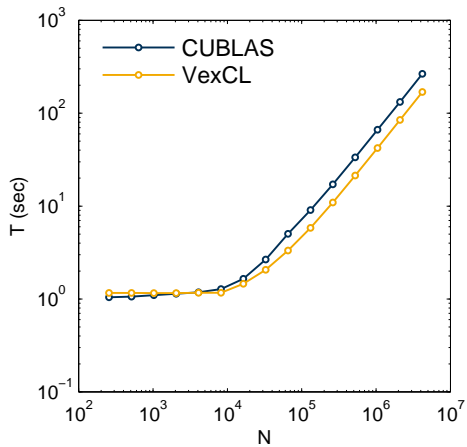
# Thrust implementation

- It is possible to fuse linear combination kernels with Thrust:

### Thrust

```
1   struct scale_sum2 {
2       const double a1, a2;
3       scale_sum2(double a1, double a2) : a1(a1), a2(a2) { }
4       template<class Tuple>
5       __host__ __device__ void operator()(Tuple t) const {
6           thrust::get<0>(t) = a1 * thrust::get<1>(t) + a2 * thrust::get<2>(t);
7       }
8   };
9
10  thrust::for_each(
11          thrust::make_zip_iterator(
12              thrust::make_tuple( x0.begin(), x1.begin(), x2.begin() )
13              ),
14          thrust::make_zip_iterator(
15              thrust::make_tuple( x0.end(), x1.end(), x2.end() )
16              ),
17          scale_sum2(a1, a2)
18          );
```

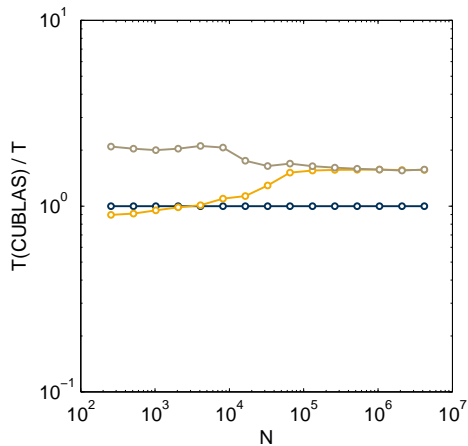### VexCL

```
1   x0 = a1 * x1 + a2 * x2;
```

- Deficiencies of naive implementation:
    - □ Runge-Kutta method uses 4 temporary state variables (here stored on GPU).
    - □ Single Runge-Kutta step results in several kernel launches.

## What if we did this manually?

- Create monolithic kernel for a single step of Runge-Kutta method.
- Launch the kernel in a loop.
- This would be 10x faster!

```
double3 lorenz_system(double r, double sigma, double b, double3 s) {
    return (double3)( sigma * (s.y - s.x),
                      r * s.x - s.y - s.x * s.z,
                      s.x * s.y - b * s.z);
}
kernel void lorenz_ensemble(
    ulong n, double dt, double sigma, double b,
    const global double *R,
    global double *X,
    global double *Y,
    global double *Z
    )
{
    for(size_t i = get_global_id(0); i < n; i += get_global_size(0)) {
        double r = R[i];
        double3 s = (double3)(X[i], Y[i], Z[i]);
        double3 k1, k2, k3, k4;

        k1 = dt * lorenz_system(r, sigma, b, s);
        k2 = dt * lorenz_system(r, sigma, b, s + 0.5 * k1);
        k3 = dt * lorenz_system(r, sigma, b, s + 0.5 * k2);
        k4 = dt * lorenz_system(r, sigma, b, s + k3);

        s += (k1 + 2 * k2 + 2 * k3 + k4) / 6;

        X[i] = s.x; Y[i] = s.y; Z[i] = s.z;
    }
}
```

# What if we did this manually?

- Create monolithic kernel for a single step of Runge-Kutta method.
- Launch the kernel in a loop.
- This would be 10x faster! But,
  - We lost odeint's generality.

```
1  double3 lorenz_system(double r, double sigma, double b, double3 s) {
2      return (double3)( sigma * (s.y − s.x),
3                        r * s.x − s.y − s.x * s.z,
4                        s.x * s.y − b * s.z);
5  }
6  kernel void lorenz_ensemble(
7      ulong n, double dt, double sigma, double b,
8      const global double *R,
9      global double *X,
10     global double *Y,
11     global double *Z
12     )
13 {
14     for(size_t i = get_global_id(0); i < n; i += get_global_size(0)) {
15         double r = R[i];
16         double3 s = (double3)(X[i], Y[i], Z[i]);
17         double3 k1, k2, k3, k4;
18
19         k1 = dt * lorenz_system(r, sigma, b, s);
20         k2 = dt * lorenz_system(r, sigma, b, s + 0.5 * k1);
21         k3 = dt * lorenz_system(r, sigma, b, s + 0.5 * k2);
22         k4 = dt * lorenz_system(r, sigma, b, s + k3);
23
24         s += (k1 + 2 * k2 + 2 * k3 + k4) / 6;
25
26         X[i] = s.x; Y[i] = s.y; Z[i] = s.z;
27     }
28 }
```

## Convert Boost.odeint stepper to a fused OpenCL kernel!

- VexCL provides $vex::symbolic<T>$ type.
- An instance of the type dumps any arithmetic operations to output stream:

```
1  vex::symbolic<double> x = 6, y = 7;
2  x = sin(x * y);
```

```
double var1 = 6;
double var2 = 7;
var1 = sin( ( var1 * var2 ) );
```

## Convert Boost.odeint stepper to a fused OpenCL kernel!

- VexCL provides vex::symbolic<T> type.
- An instance of the type dumps any arithmetic operations to output stream:

```
1  vex::symbolic<double> x = 6, y = 7;
2  x = sin(x * y);
```

```
double var1 = 6;
double var2 = 7;
var1 = sin( ( var1 * var2 ) );
```

- The idea is very simple:
  - Record sequence of arithmetic expressions of an algorithm.
  - Generate OpenCL kernel from the captured sequence.

# Record operations performed by Boost.odeint stepper

### 1. State type

```
1  typedef vex::symbolic< double > sym_vector;
2  typedef std::array<sym_vector, 3> sym_state;
```

### 2. System functor

```
3  struct lorenz_system {
4      const sym_vector &R;
5      lorenz_system(const sym_vector &R) : R(R) {}
6
7      void operator()(const sym_state &x, sym_state &dxdt, double t) const {
8          dxdt[0] = sigma * (x[1] − x[0]);
9          dxdt[1] = R * x[0] − x[1] − x[0] * x[2];
10         dxdt[2] = x[0] * x[1] − b * x[2];
11     }
12 };
```

# Record operations performed by Boost.odeint stepper

### 3. Stepper

```
13   odeint :: runge_kutta4<
14           sym_state /*state*/,        double /*value*/,
15           sym_state /*derivative*/, double /*time*/,
16           odeint :: range_algebra, odeint :: default_operations
17           > stepper;
```

### 4. Record one step of Runge-Kutta method

```
18   std :: ostringstream lorenz_body;
19   vex :: generator :: set_recorder (lorenz_body);
20
21   sym_state sym_S = {{ sym_vector(sym_vector::VectorParameter),
22                        sym_vector(sym_vector::VectorParameter),
23                        sym_vector(sym_vector::VectorParameter) }};
24   sym_vector sym_R(sym_vector::VectorParameter, sym_vector::Const);
25
26   lorenz_system sys(sym_R);
27   stepper.do_step(std :: ref(sys), sym_S, 0, dt);
```

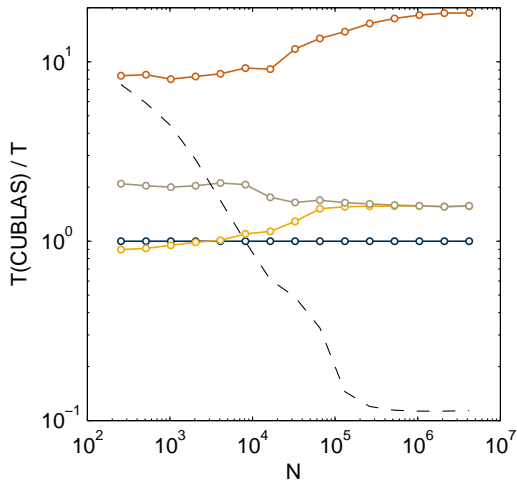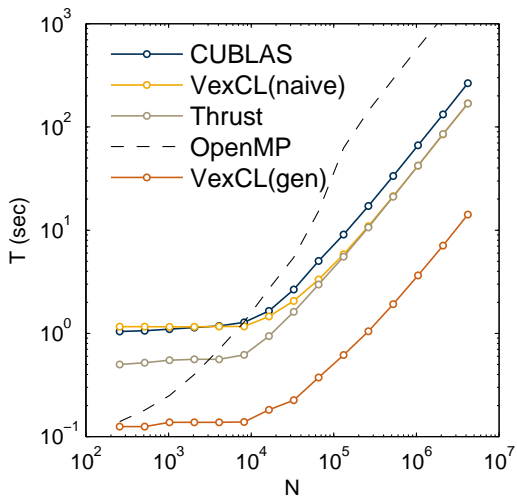### 5. Generate and use OpenCL kernel

```
28  auto lorenz_kernel = vex::generator::build_kernel(ctx, "lorenz", lorenz_body.str(),
29        sym_S[0], sym_S[1], sym_S[2], sym_R);
30
31  vex::vector<double> X(ctx, n), Y(ctx, n), Z(ctx, n), R(ctx, n);
32
33  X = Y = Z = 10;
34  R = Rmin + (Rmax − Rmin) * vex::element_index() / (n − 1);
35
36  for(double t = 0; t < t_max; t += dt) lorenz_kernel(X, Y, Z, R);
```

# The restrictions

- Algorithms have to be embarrassingly parallel.
- Only linear flow is allowed (no conditionals or data-dependent loops).
- Some precision may be lost when converting constants to strings.

# Performance of the generated kernel

# Projects using VexCL

AMGCL — algebraic multigrid implementation:

- https://github.com/ddemidov/amgcl

Antioch — A New Templated Implementation Of Chemistry for Hydrodynamics:

- https://github.com/libantioch/antioch

Boost.odeint — numerical solution of Ordinary Differential Equations:

- http://odeint.com

# Summary

- VexCL allows to write compact and readable code without sacrificing performance.
- Its code generator allows to convert generic C++ code to OpenCL at runtime:
  - □ Reduces global memory I/O
  - □ Reduces number of kernel launches
  - □ Facilitates code reuse

- Supported compilers (don't forget to enable C++11 features):
  - □ GCC v4.6
  - □ Clang v3.1
  - □ MS Visual C++ 2010

- https://github.com/ddemidov/vexcl