# Thesis contributions

- APPy: Annotated Parallelism for Python on GPUs
  - [CC24] Parallelize Python loops and tensor expressions on GPUs
- ReACT: Redundancy-Aware Code Generation for Tensor Expressions
  - [PACT22] Redundancy elimination when fusing sparse/dense tensor operators
- **Intrepydd: Performance, Productivity, and Portability for Data Science Application Kernels**
  - [Onward!20] Compile Python/NumPy to C++ with high-level optimizations

Georgia Tech

# Problem statement: desired input and output

- Desired input: whole kernel in Python (control flow is fine)

- Desired output: C++ code

1. it = 0
2. **while** it < max_iter:
3.     u = 1.0 / x
4.     v = c * (1 / (K.T @ u))
5.     x = ((1 / r) * K) @ v
6.     it += 1

```cpp
222  py::array_t<double> train(py::array_t<int> A, py::array_t<double> F,
223                             int iterations) {
224      /* Declarations */
225      double *F_p_data_ptr_pydd;
226      double *grad_data_ptr_pydd;
227      int64_t N;
228      int n;
229      int person;
230      py::array_t<double> grad;
231      py::array_t<double> F_p;
232      double ll;
233      int __var7;
234
235      N = pydd::shape(A, 0);
236      for (int _i = 0; _i < iterations; _i += 1) {
237          n = _i;
238          for (int _i = 0; _i < N; _i += 1) {
239              person = _i;
240              grad = gradient(F, A, person);
241              F_p = pydd::get_row(F, person);
242              pydd::compatibility_check(F_p, grad);
243              F_p_data_ptr_pydd = F_p.mutable_data();
244              int F_p_shape0 = pydd::shape(F_p, 0);
245              // int F_p_shape0 = pydd::shape(F, 1);
246              // F_p_data_ptr_pydd = (double*)F.mutable_data() + person*F_p_shape0;
247
248              grad_data_ptr_pydd = grad.mutable_data();
249              for (int _i = 0; _i < F_p_shape0; _i += 1) {
250                  __var7 = _i;
251                  pydd::setitem_1d(
252                      F_p_data_ptr_pydd,
253                      (pydd::getitem_1d(F_p_data_ptr_pydd, __var7) +
254                      (0.005 * pydd::getitem_1d(grad_data_ptr_pydd, __var7))),
255                      __var7);
256              };
257
258
259              pydd::set_row(F, person, pydd::maximum(0.001, F_p));
260
261          };
262          ll = log_likelihood(F, A);
263      };
264      return F;
265  }
```

Georgia Tech

# Compilation Pipeline: From Intrepydd to C++

**Intrepydd source code**

```
1.    def foo(xs: Array(double, 2)) -> double:
             ...
2.        for i in range(shape(xs, 0)):
3.          for j in range(shape(xs, 1)):
4.             sum += xs[i, j]
5.                                        ...
```

Georgia Tech

# Compilation Pipeline: From Intrepydd to C++

**Intrepydd source code**

```
1.    def foo(xs: Array(double, 2)) -> double:
            ...
2.      for i in range(shape(xs, 0)):
3.        for j in range(shape(xs, 1)):
4.          sum += xs[i, j]
5.                                  ...
```

**Intrepydd compiler**

**Resulting C++ code**

```
1.    Array<double>* foo(Array<double>* xs) {
2.    ...
3.    for (int i = 0; i < pydd::shape(xs, 0); i += 1) {
4.      for (int j = 0; j < pydd::shape(xs, 1); j += 1) {
5.        sum += xs.data()[i*pydd::shape(xs, 1)+j];
6.          ...
```

Georgia Tech.

# Code Optimization

- High-level Optimizations in AOT compilation
  - Loop invariant code motion (LICM OPT)
  - Dense & Sparse Array Operator Fusion (Array OPT)
  - Array allocation and slicing optimization (Memory OPT)

Georgia Tech.

# Code Optimization: LICM

c: sparse
K, u: dense

```
1.    it = 0
2.    while it < max_iter:
3.        u = 1.0 / x
4.        v = c * (1 / (K.T @ u))   # SDDMM
5.        x = ((1 / r) * K) @ v
6.        it += 1
```

**Intrepydd source code (Sinkhorn)**

```
1.    it = 0
2.    # Hoisted loop-invariant expressions
3.    tmp1 = K.T
4.    tmp2 = (1 / r) * K
5.    while it < max_iter:
6.        u = 1.0 / x

7.        v = empty_like(c)
8.        # Fused loop iterating over non-zero elements
9.        for row, col, val in c.nonzero_elements():
10.           tmp3 = 0.0
11.           for idx in range(shape(tmp1, 1)):
12.               tmp3 += tmp1[row, idx] * u[idx, col]
13.           tmp4 = val * (1 / tmp3)
14.           spm_set_item(v, tmp4, row, col)

15.       x = spmm_dense(tmp2, v)

16.       it += 1
```

**Transformed code**

Georgia Tech

# Code Optimization: Sparse Operator Fusion

c: sparse
K, u: dense

```
1.    it = 0
2.    while it < max_iter:
3.        u = 1.0 / x
4.        v = c * (1 / (K.T @ u)) # SDDMM
5.        x = ((1 / r) * K) @ v
6.        it += 1
```

SDDMM: masked matmul

**Intrepydd source code (Sinkhorn)**

```
1.    it = 0
2.    # Hoisted loop-invariant expressions
3.    tmp1 = K.T
4.    tmp2 = (1 / r) * K
5.    while it < max_iter:
6.        u = 1.0 / x
7.        v = empty_like(c)
8.        # Fused loop iterating over non-zero elements
9.        for row, col, val in c.nonzero_elements():
10.           tmp3 = 0.0
11.           for idx in range(shape(tmp1, 1)):
12.               tmp3 += tmp1[row, idx] * u[idx, col]
13.           tmp4 = val * (1 / tmp3)
14.           spm_set_item(v, tmp4, row, col)

15.       x = spmm_dense(tmp2, v)

16.       it += 1
```

**Transformed code**

Georgia Tech

# Code Optimization: Dense Operator Fusion

c: sparse
K, u: dense

```
1.    it = 0
2.    while it < max_iter:
3.        u = 1.0 / x
4.        v = c * (1 / (K.T @ u))   # SDDMM
5.        x = ((1 / r) * K) @ v
6.        it += 1
```

SDDMM: masked matmul

**Intrepydd source code (Sinkhorn)**

```
1.    it = 0
2.    # Hoisted loop-invariant expressions
3.    tmp1 = K.T
4.    tmp2 = (1 / r) * K
5.    while it < max_iter:
6.        u = 1.0 / x
7.        v = empty_like(c)
8.        # Fused loop iterating over non-zero elements
9.        for row, col, val in c.nonzero_elements():
10.           tmp3 = 0.0
11.           for idx in range(shape(tmp1, 1)):
12.               tmp3 += tmp1[row, idx] * u[idx, col]
13.           tmp4 = val * (1 / tmp3)
14.           spm_set_item(v, tmp4, row, col)

15.       x = spmm_dense(tmp2, v)

16.       it += 1
```

**Transformed code**

65

Georgia Tech

# Experimental Methodology

## Benchmark Applications

- A subset of Python based data analytics applications from a recent DARPA program
- Mix of non-library call and library call dominated applications
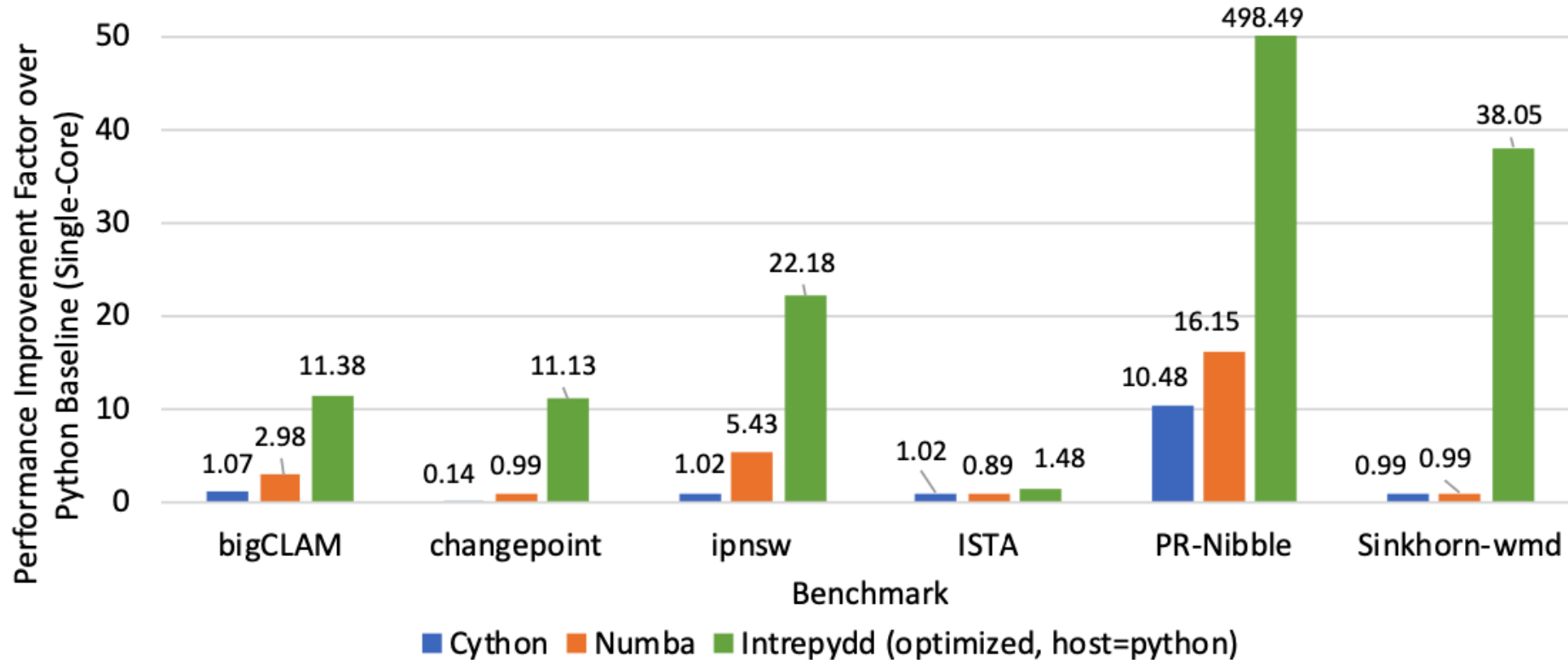
## Test machine

- Dual Intel Xeon Silver 4114 CPU @ 2.2GHz with 192GB of main memory and hyperthreading disabled

## Comparisons

- Baseline idiomatic Python 3.7.6
- Cython
- Numba

Georgia Tech.

# Intrepydd Sequential Performance



**Intrepydd offers 20.7x speedup on average (geomean) over baseline Python**

Georgia Tech

# Code Optimization

- High-level Optimizations in AOT compilation
    - Loop invariant code motion (LICM OPT)
    - Dense & Sparse Array Operator Fusion (Array OPT)
    - Array allocation and slicing optimization (Memory OPT)

- Impact on performance by each OPT

| | Primary Kernel execution times (seconds) | | | |
|---|---|---|---|---|
| Benchmark | Intrepydd | Intrepydd (+LICM OPT) | Intrepydd (+Array OPT) | Intrepydd (+Memory OPT) |
| bigCLAM | 2.558 | 2.557 | 1.541 | 1.086 |
| changepoint | 1.472 | 1.469 | 1.466 | 1.471 |
| ipnsw | 1.679 | 0.786 | 0.786 | 0.786 |
| ISTA | 79.362 | 18.732 | 18.473 | 18.509 |
| PR-Nibble | 0.831 | 0.114 | 0.106 | 0.106 |
| sinkhorn-wmd | 47.612 | 47.395 | 1.225 | 1.220 |

Georgia Tech

# Intrepydd summary

- We present Intrepydd, a Python-based programming system, which is designed to enable data scientists to write application kernels with high performance, productivity, and portability

- We implement a number of high-level compiler optimizations during the compilation

- We evaluate the performance of Intrepydd using 6 data science kernels and show significant single-core performance improvements of Intrepydd relative to vanilla Python/NumPy (1.5× to 498.5×), Cython (1.5× to 47.5×) and Numba (1.7× to 38.1×)

Georgia Tech