# The Waverun Programming Model for FPGA-Based Heterogeneous Architectures

## Abhi D.R. and Ron Sass

{adevalap,rsass}@uncc.edu

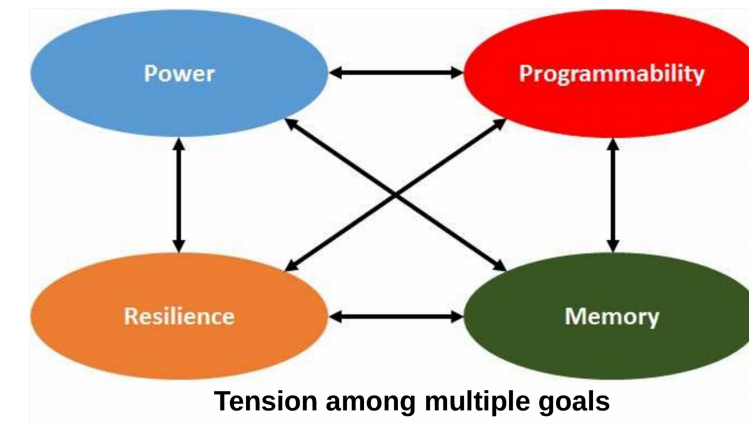Reconfigurable Computing Systems Lab: http://www.rcs.uncc.edu

## Motivation

The PyDAC runtime for a FPGA-based heterogeneous architecture demonstrated that algorithms using the Divide-and-Conquer design pattern can balance multiple --- sometimes conflicting --- system goals.



Tension among multiple goals

### Key Question:

*Can the same success be extended to another class of applications that use, for example, the Wavefront design pattern?*
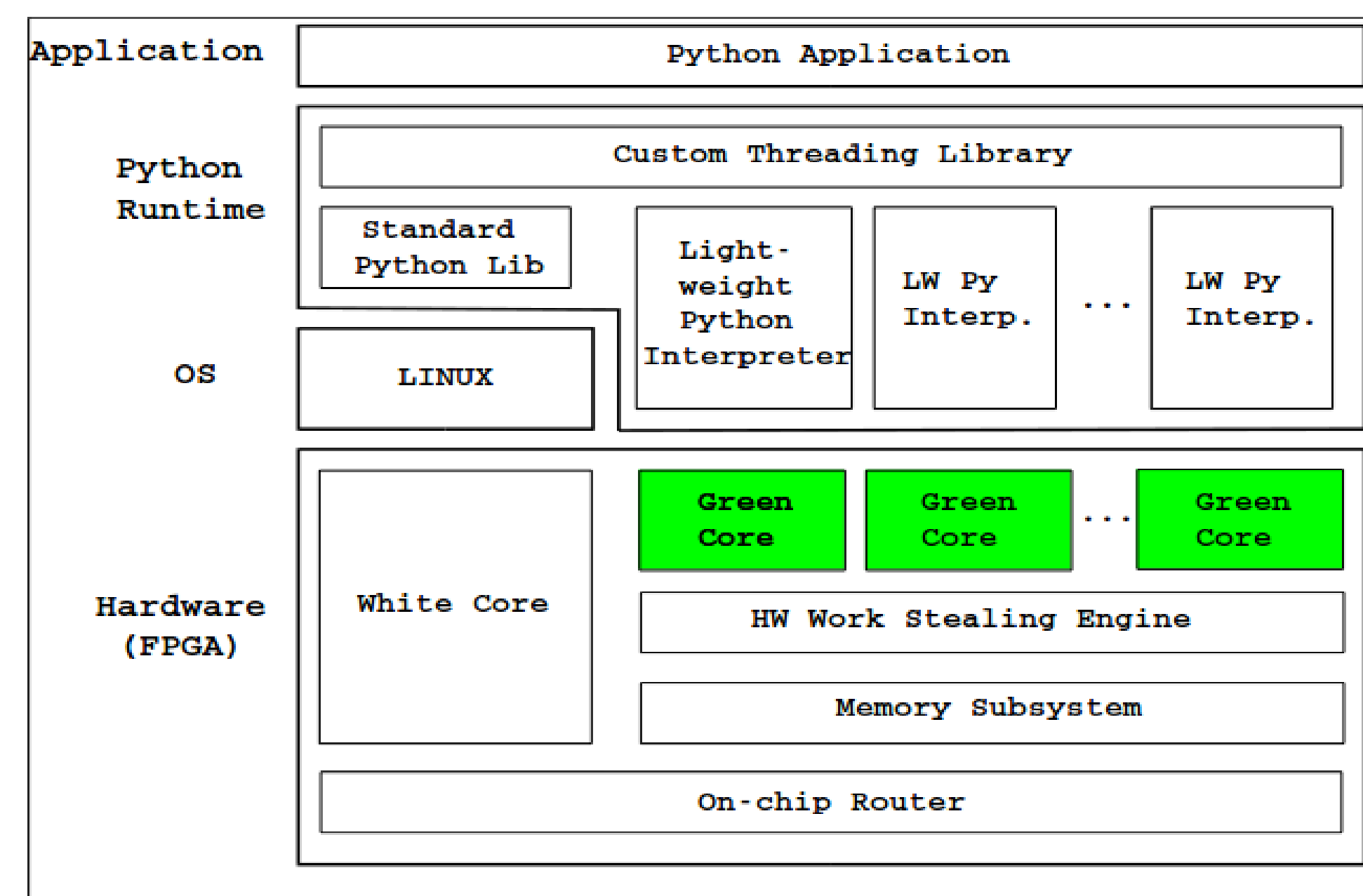
To test this, we left the hardware unchanged and investigated four well-known wavefront algorithms: LU decomposition, QR decomposition, Cholesky decomposition and matrix-matrix multiplication.

Success is defined as:

♦ A simple application/runtime interface
♦ Efficient use of the memory subsystem

(The hardware does not change so power and resilience are not impacted by these tests.)

### FPGA Heterogeneous Architecture



♦ `Match` and `Post` method for memory access
♦ Object store and no linear array of memory
Currently implemented in software but part of heterogeneous hardware

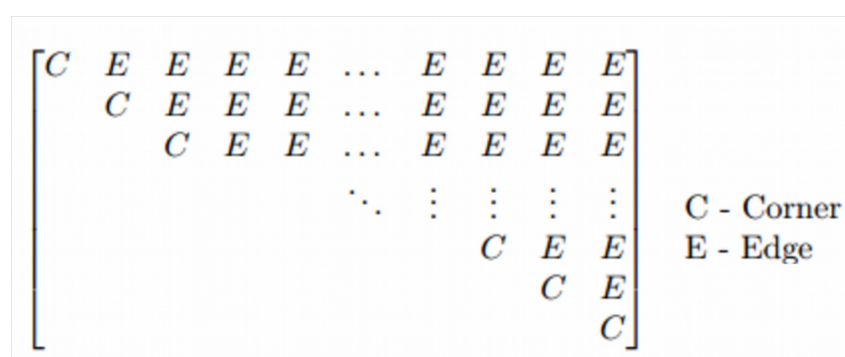## Examples

```
1   class dpbtrf(waverun):
2       def com_corner(self,Ainput,Sdatajk,Sdatakk):
3           return Ainput-sum([(p**2*q) for p,q in zip(Sdatajk,Sdatakk)])
4       def com_edge(self,Ainput,Sdatajk,Sdataik,chofact):
5           return (Ainput-sum([(p**2*q*r) for p,q,r in zip(Sdatajk,Sdatakk,\
6                   Sdataik)]))/chofact
7
8       def Corner(self,index,step):
9           Res=self.com_corner(Ainput,jkdata,kkdata)#Compute
10          self.dptrb_post(index,step,Res)#Post data
11          return
12      def Edge(self,index,step):
13          Res=self.com_edge(Ainput,jkdata,kkdata,ikdata,chofact)#Compute
14          self.dptrb_post(index,step,Res)#Post data
15          return
16
17      def getFunc(self):
18          Func={'00':self.Corner,'0*':self.Edge}
19          return Func
20      def primer(self):
21          for i in range (self.data.shape[0]/block):
22              for j in range (self.data.shape[1]/block):
23                  if(i<=j):
24                      self.post(self.data[i][j],'Ainput',[i+1,j+1],0)
25                  if(i==0):
26                      self.post(self.data.shape[0]-j-1,'maxi',[i+1,j+1],0)
27                  else:
28                      self.post(-1,'maxi',[i+1,j+1],0)
29                      self.post(0,'Result',[i+1,j+1],0)
```

Figure 2: Cholesky Decomposition

$$\begin{bmatrix} C & E & E & E & \dots & E & E & E \\ & C & E & E & \dots & E & E & E \\ & & C & E & E & \dots & E & E \\ & & & \ddots & \vdots & \vdots & \vdots & \vdots \\ & & & & C & E & E \\ & & & & & C & E \\ & & & & & & C \end{bmatrix}$$

C - Corner
E - Edge

**Computation:**

• Corner
$$D_{jj} = A_{jj} - \sum_{k=1}^{j-1} U_{jk}^2 D_{kk}$$

• Edge
$$U_{ij} = A_{ij} - \sum_{k=1}^{j-1} U_{ik} U_{jk} D_{kk}$$

```
1   class dgemm(waverun):
2       def Mul(self,Ainput,Binput,Cold):
3           return np.add(np.dot(Ainput,Binput),Cold)
4
5       def Corner(self,index,step):
6           #Match data
7           Cnew=self.Mul(Ainput,Binput,Cold)
8           #Post results
9           return
10      def Xedge(self,index,step):
11          #Match data
12          Cnew=self.Mul(Ainput,Binput,Cold)
13          #Post results
14          return
15      def Yedge(self,index,step):
16          #Match data
17          Cnew=self.Mul(Ainput,Binput,Cold)
18          #Post results
19          return
20      def Interior(self,index,step):
21          #Match data
22          Cnew=self.Mul(Ainput,Binput,Cold)
23          #Post results
24          return
25
26      def getFunc(self):
27          Func={'00':self.Corner,'*0':self.Xedge,'0*':self.Yedge,'**':self.Interior}
28          return Func
29
30      def primer(self):
31          for i in range (self.data[0].shape[0]/block):
32              for j in range (self.data[1].shape[1]/block):
33                  #Post input data
```

$$\begin{bmatrix} C & X & X & X & \dots & X & X & X \\ Y & I & I & I & \dots & I & I & I \\ Y & I & I & I & \dots & I & I & I \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ Y & I & I & I & \dots & I & I & I \\ Y & I & I & I & \dots & I & I & I \\ Y & I & I & I & \dots & I & I & I \end{bmatrix}$$

C - Corner
X - Xedge
Y - Yedge
I - Interior

Figure 1: Matrix Multiplication

## Runtime

♦ Provides simple user interface

**Schedule** and **Gather** methods to schedule the task and gather the result

♦ One program for different types of cores
♦ Manages memory keys

**Programmer responsibility:**
♦ Simple object-oriented interface to runtime
♦ Provide function based on the index
♦ Block of input data via Python method

## Results

A simple application/runtime interface with two user methods is implemented

```
Cholesky Decomposition
1   Amat=np.random.randint(num,size=(row,column))
2   myapp=dpbtrf(Amat,block)
3   myapp.schedule(row/block)
4   myapp.gather('Result')
```
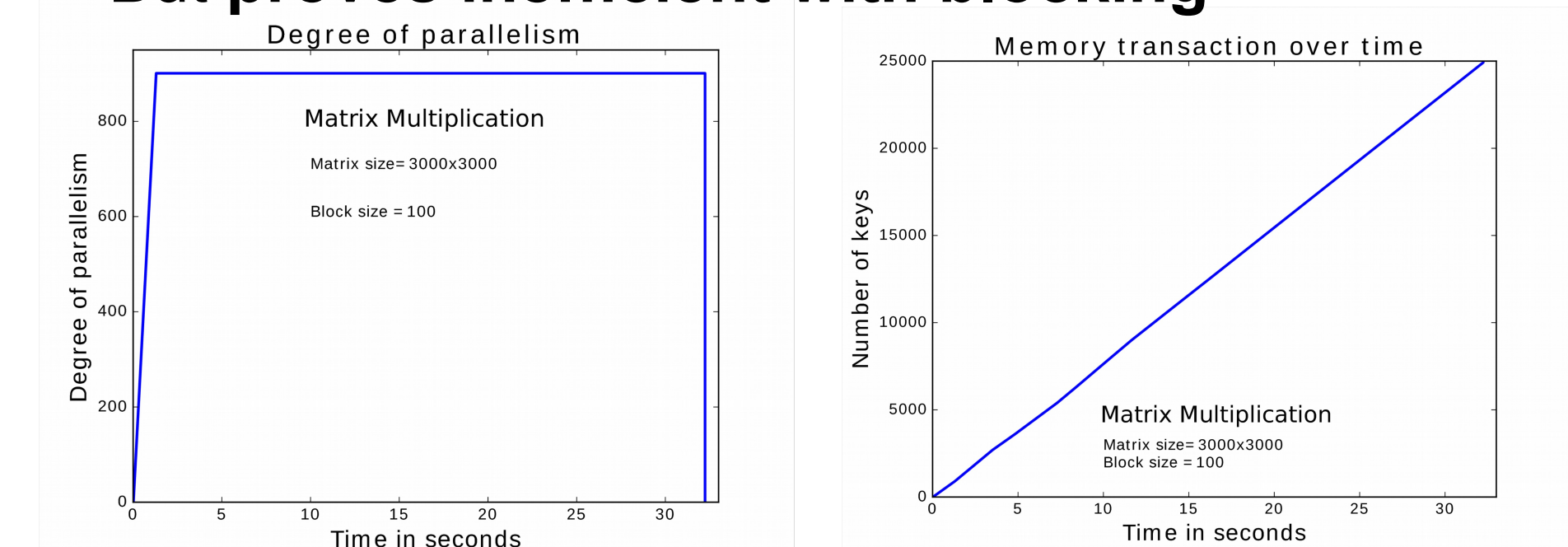
```
Matrix Multiplication
1   Amat=np.random.randint(num,size=(row,column))
2   Bmat=np.random.randint(num,size=(row,column))
3   myapp=dgemm((Amat,Bmat),block)
4   myapp.schedule(row/block)
5   myapp.gather('Result')
```

**User application using the runtime interface**

♦ Cholesky, LU, QR, and Matrix multiplication algorithms were implemented using "wavefront" design pattern
♦ Cholesky, LU, and QR works well with point-point computation
♦ But **proves inefficient with blocking**



**Results of matrix multiplication with blocking**

Hardware memory subsystem is tested for transaction bandwidth of range 20k-40k. The runtime efficiently uses the system by generating keys of similar range.

## Background

**Pydac previously implemented:**
♦ Divide and conquer
♦ Proven to be efficient for certain algorithms

B. Huang, R. Sass, N. Debardeleben and S. Blanchard, "Harnessing Unreliable Cores in Heterogeneous Architecture: The PyDac Programming Model and Runtime," *2014 44th Annual IEEE/IFIP International Conference on Dependable Systems and Networks*, Atlanta, GA, 2014, pp. 744-749.

**Design Patterns**
Intel TBB (Parallel patterns)
♦ Divide and Conquer
♦ Wavefront
♦ Agglomeration
♦ Elementwise
♦ Reduction

## Future Work

♦ Developing graph related algorithms for waverun
♦ Testing the runtime with the heterogeneous hardware