# Hardware, tool, and algorithm facets in next generation low power HPEC (High Performance Embedded Computing) that should shape benchmarking and evaluation strategy

## Disclaimer

These are research projects, not all in Reservoir products.

No warranties.

Patent pending technologies.

# Game Over or Next Level?

"A high-level performance-portable programming model is the only way to restart the virtuous cycle..."



THE FUTURE OF
COMPUTING PERFORMANCE

Game Over or
Next Level?

NATIONAL RESEARCH COUNCIL
OF THE NATIONAL ACADEMIES

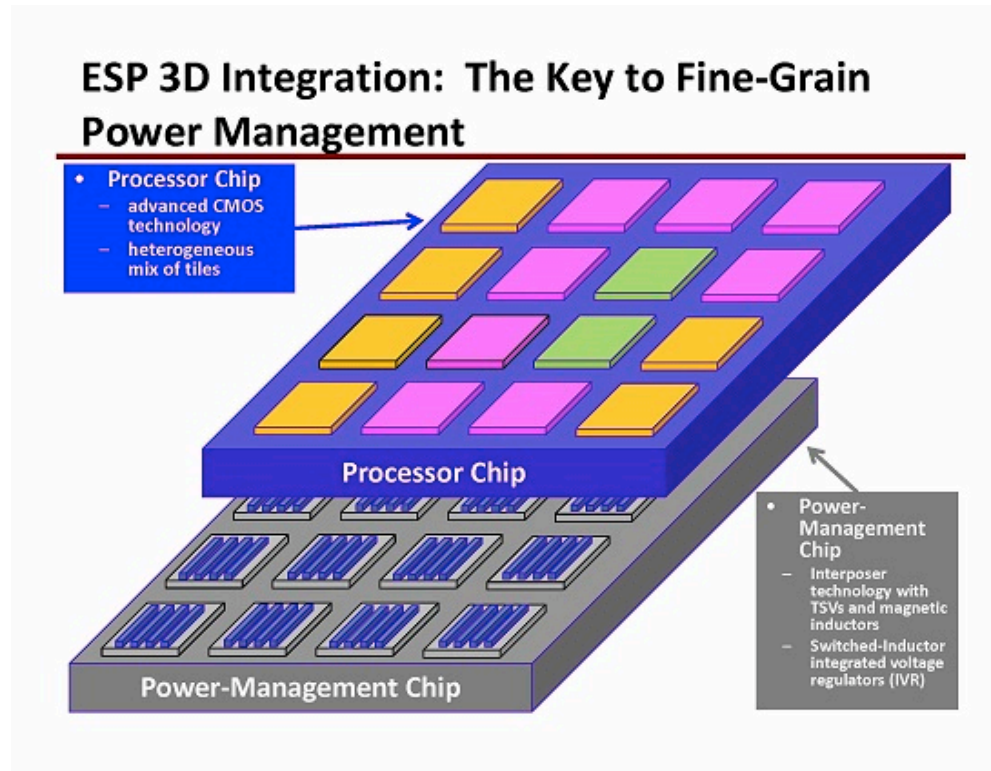## Outline

Hardware Facets

Tool Facets

Algorithm Facets

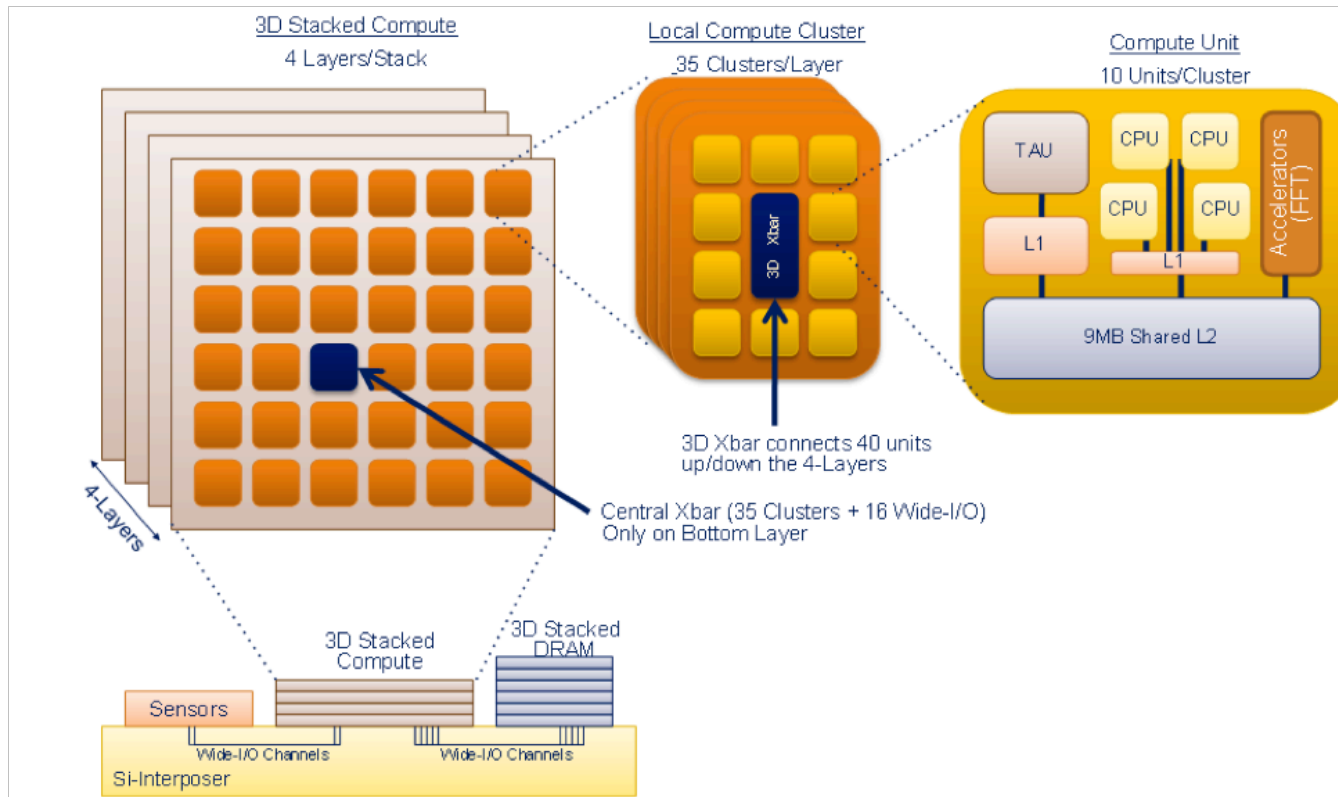Benchmarking and Evaluation Implications

# Hardware Facets

# Nimble, Fine Grained Per Core Power Controls



"Embedded Scalable Processor (ESP)"
Columbia University
Luca Carloni (PI), 2014

# More Cores, Deeper Hierarchy, NTV



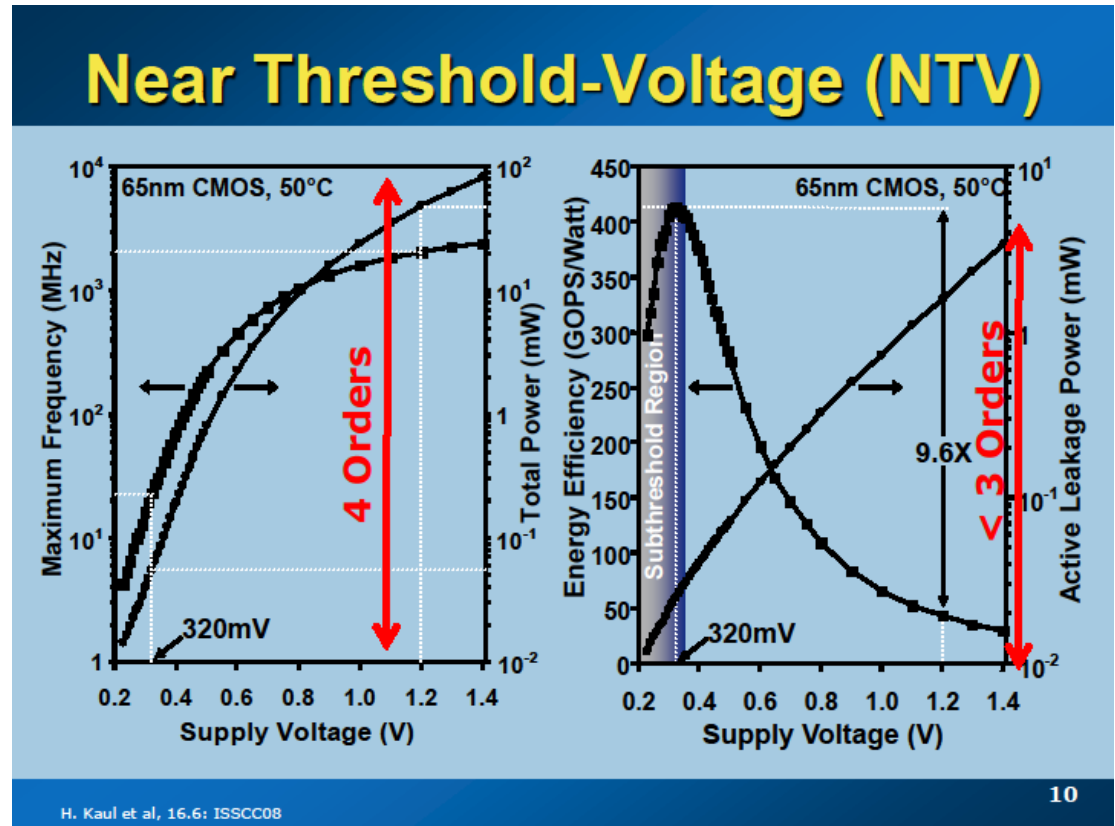"Energy Efficient 3D Near-Threshold Computing Systems for Future
Embedded Applications"
Trevor Mudge
University of Michigan, 2014

# NTV Operation

## NTV Benefit:

4X Increase in Power Efficiency

## NTV Cost:

Need more CONCURENCY
Greater leakage variation
Greater f_max variation
Bigger relative register files



Shekhar Borkar
IPDPS keynote
2013

# NTV ➜ Cores running at varied clocks

Cost:

More complex mapping
and LOAD BALANCING



Shekhar Borkar
IPDPS keynote
2013

# Energy Cost Trends: Arithmetic vs. Communications

Arithmetic

Communications



Log Scale (improving exponentially;
i.e., Moore's Law)

Linear Scale (improving only linearly!)

Needed:
Communication-Avoiding Compilation,
Communication-Avoiding Algorithms

# Other Hardware Trends

Heterogeneity

Compute Everywhere (in Network, near Memory, ...)

Sea of Tiles

Specialization / Dark Silicon

Explicit Communication (DMA, rDMA, ...)

Explicit Scratchpad hierarchy management

...

# Tool Facets

# R-Stream : More CONCURRENCY

# Targeting Task Graph Execution Models

Move beyond barrier synchronization

- Expose additional CONCURRENCY
- Increase opportunities for LOAD BALANCING
- Support runtimes that scale to 1000's of cores

Example: Givens QR mapping only requires neighbor syncs

unnecessary

# Benefits

Better efficiency, shorter latency, lower power usage

Dynamic optimizations (locality, power, energy proportionality)



2 threads, 25 tasks

# Open Community Runtime



https://01.org/projects/open-community-runtime

# Polyhedral Mapper in R-Stream

# Communication-Avoiding Scheduling

Affine Scheduling – JPLCVD

- Jointly optimize parallelism, locality, contiguity, vectorization and data layout

Local Memory Buffer Management

- Scratchpad and/or "Virtual Scratchpad"

Minimize copying and memory traffic in local memory

- "Right size" local buffers and round-robin reuse

# Communication-Avoiding Compilation

```
for k=1,n {
    for j=1,n {
        for i=1,n {
            c[3*(i-1)+1] = A[k][j][i+1];
            c[3*(i-1)+2] = A[k][j][i+2];
            c[3*(i-1)+3] = A[k][j][i+3];
                ...
            X[k][j][i] = c[3*(i-1)+1] *
                         c[3*(i-1)+2] +
                         c[3*(i-1)+3];
        }
    }
}
```
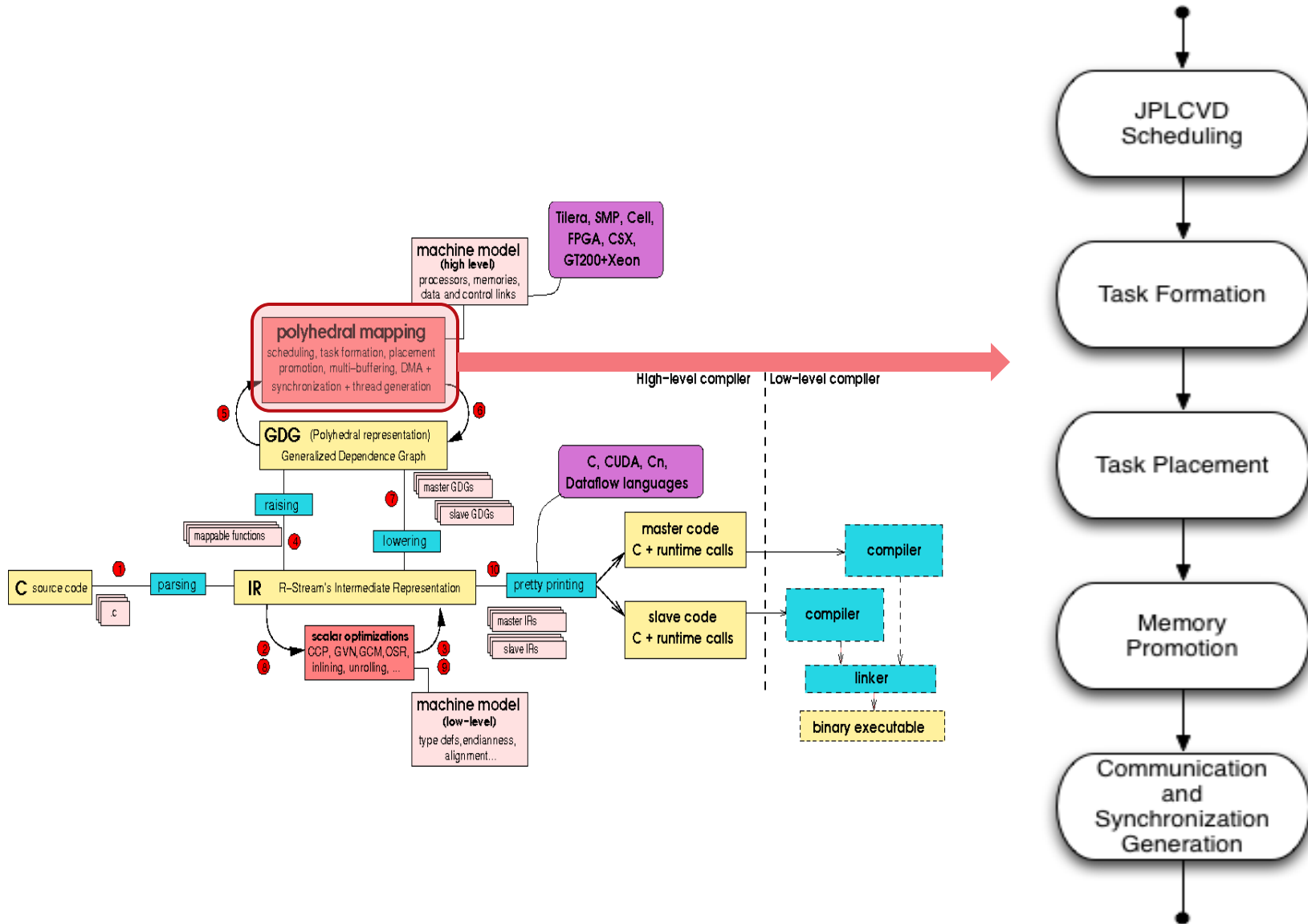
Original Code (Covariance Calculation from STAP)

# Communication-Avoiding Compilation

```
for k=1,n {
   for j=1,n {
      for i=1,n {
         if (i == 1) {
            c[3*(i-1)+1] = A[k][j][i+1];
            c[3*(i-1)+2] = A[k][j][i+2];   // S1
            c[3*(i-1)+3] = A[k][j][i+3];   // S2
               ...
            X[k][j][i] = c[3*(i-1)+1] *
                         c[3*(i-1)+2] +
                         c[3*(i-1)+3];
         }
         if (i == 2) {
            c[3*(i-1)+1] = A[k][j][i+1];   // S3
            c[3*(i-1)+2] = A[k][j][i+2];   // S4
            c[3*(i-1)+3] = A[k][j][i+3];   // S5
               ...
            X[k][j][i] = c[3*(i-1)+1] *
                         c[3*(i-1)+2] +
                         c[3*(i-1)+3];    // S6
         }
         if (i > 2) {
            ...
         }
      }
   }
}
```
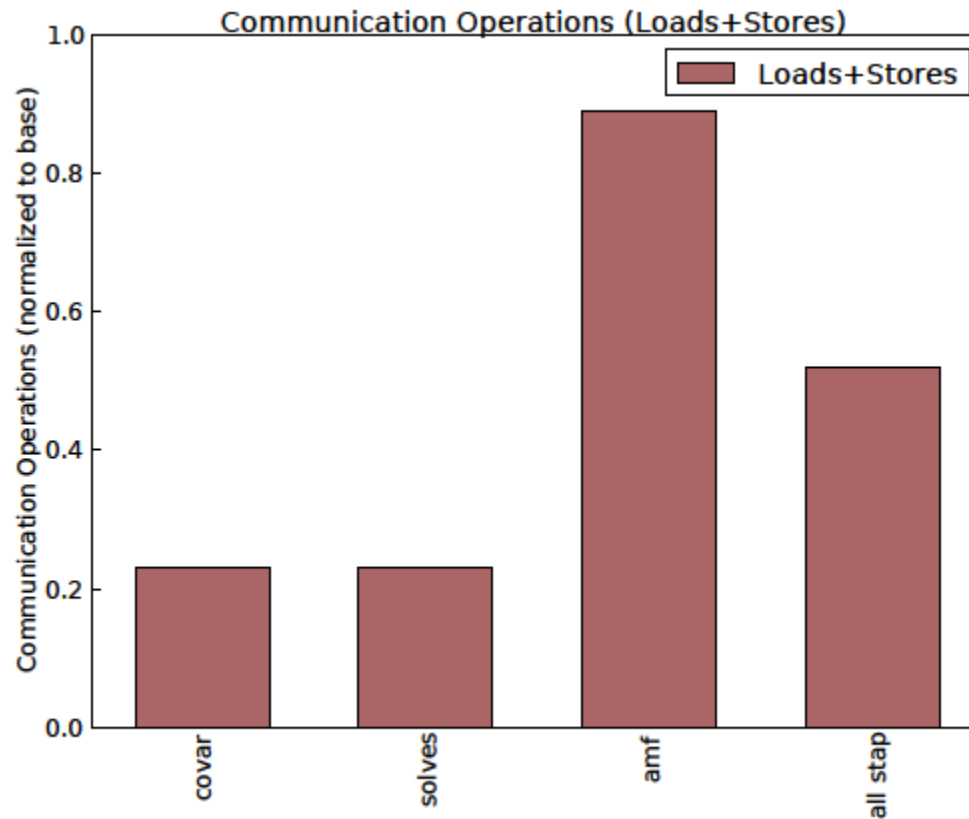
Code with two iterations of i loop peeled for illustration

# Communication-Avoiding Compilation

```
for k=1,n {
   for j=1,n {
      for i=1,n {
         if (i == 1) {
            c[3*(i-1)+1] = A[k][j][i+1];
            c[3*(i-1)+2] = A[k][j][i+2];  // S1
            c[3*(i-1)+3] = A[k][j][i+3];  // S2
                  ...
            X[k][j][i] = c[3*(i-1)+1] *
                         c[3*(i-1)+2] +
                         c[3*(i-1)+3];
         }
         if (i == 2) {
            //c[3*(i-1)+1] = A[k][j][i+1];  // S3 : Removed - same as S1 read (S3->S1 RAR)
            //c[3*(i-1)+2] = A[k][j][i+2];  // S4 : Removed - same as S2 read (S4->S2 RAR)
            c[(i-1)\%3] = A[k][j][i+3];     // S5: Access fn change: c[3*(i-1)+3] -> c[((3*(i-1)+3)/3-1)\%3]
                  ...
            X[k][j][i] = c[3*(i-2)+2] *  // S6: Changed access fn to point to element of c loaded in S1
                                         //   c[3*(i-1)+1] -> c[3*(i-2)+2] (S6->S3 RAW + S3->S1 RAR)
                         c[3*(i-2)+3] +  // S6: Changed access fn to point to element of c loaded in S2
                                         //   c[3*(i-1)+2] -> c[3*(i-2)+3] (S6->S4 RAW + S4->S2 RAR)
                         c[(i-1)\%3];    // S6: Access fn change: c[3*(i-1)+3] -> c[((3*(i-1)+3)/3-1)\%3]
         }
         if (i > 2) {
            ...
         }
      }
   }
}
```

Code with communication–avoiding optimizations applied

---

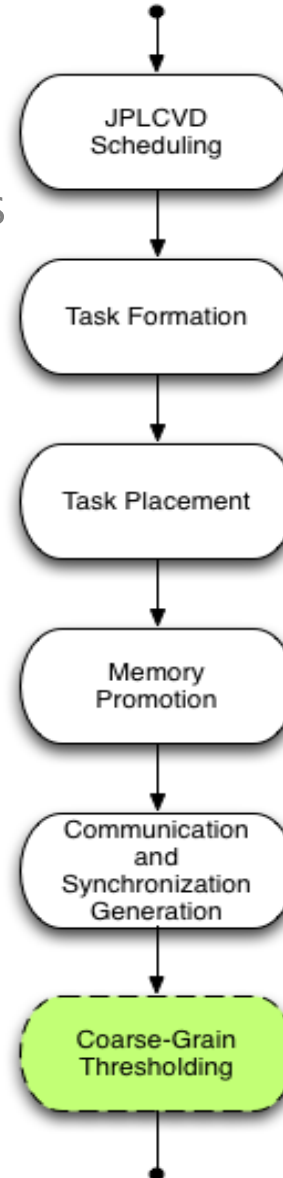# Communications-Savings (GPGPU – Kayla)



New R-Stream Optimizations Decrease Communications by 2X

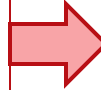# Energy Proportional Scheduling in R-Stream

## Coarse-grain Thresholding

● Generate lightweight runtime checks to apply appropriate level of energy optimization on key system components using the Power API
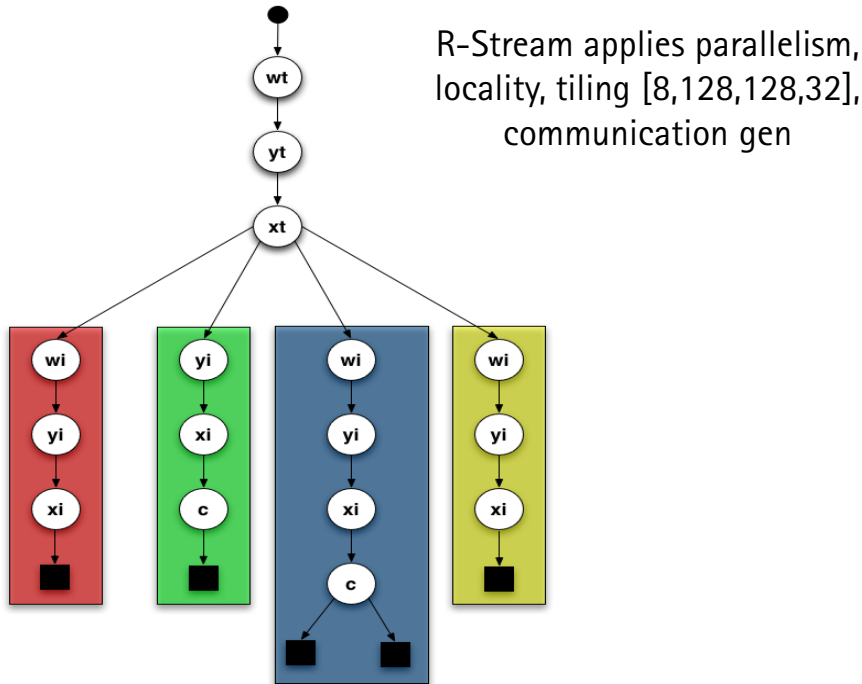
# EPS in R-Stream

```
for(w=0;w<W;w++)      // W=56
 for(y=0;y<Y0;y++)    // Y0=1021
  for(x=0;x<X0;x++)   // X0=1021
   for(c=0;c<C;c++)   // C=32
    out[w][y][x] +=
    in[y][x][c]*filt[w][0][0][c]+
    in[y+1][x][c]*filt[w[[1][0][c]+
       …
    in[y+2][x+2][c]*filt[w][2][2][c];
```

R-Stream applies parallelism, locality, tiling [8,128,128,32], communication gen



Betatree representation in R-Stream

```
for(wt=0;wt<=6;wt++){
 for(yt=0;yt<=7;yt++){
  for(xt=0;xt<=7;xt++){
   for(wi=…){
    for(yi=0…){
     for(xi=0…){
      … /* load array out */
}}}
   for(yi=…){
    for(xi=0…){
     for(c=0…){
      … /* load array in */
}}}
   for(wi=0;wi<=7;wi++){
    for(yi=0;yi<=min(…,127);yi++){
     for(xi=0;xi<=min(…,127);xi++){
      for(c=0;c<=31;c++){
       … /* load array filt */
       out_l[8*wt+wi][…][…] += …
}}}}
   for(wi=…){
    for(yi=0…){
     for(xi=0…){
      … /* store array out */
}}}
```
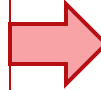
# EPS in R-Stream

```
for(w=0;w<W;w++)        // W=56
 for(y=0;y<Y0;y++)     // Y0=1021
  for(x=0;x<X0;x++)    // X0=1021
   for(c=0;c<C;c++)    // C=32
    out[w][y][x] +=
    in[y][x][c]*filt[w][0][0][c]+
    in[y+1][x][c]*filt[w[[1][0][c]+
        …
    in[y+2][x+2][c]*filt[w][2][2][c];
```
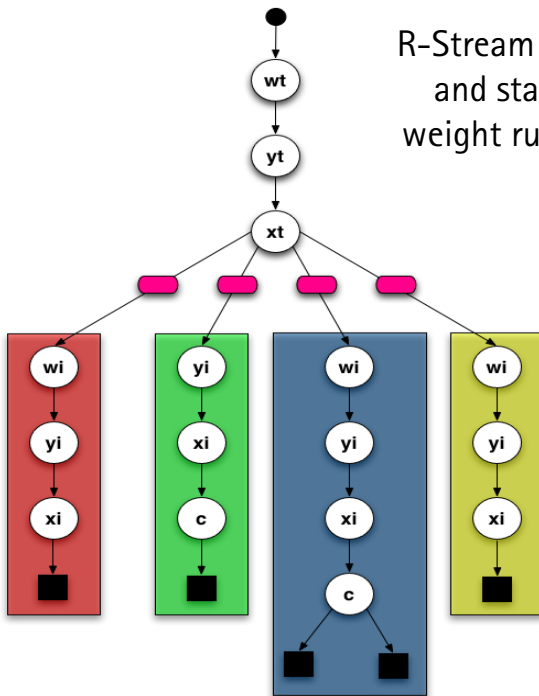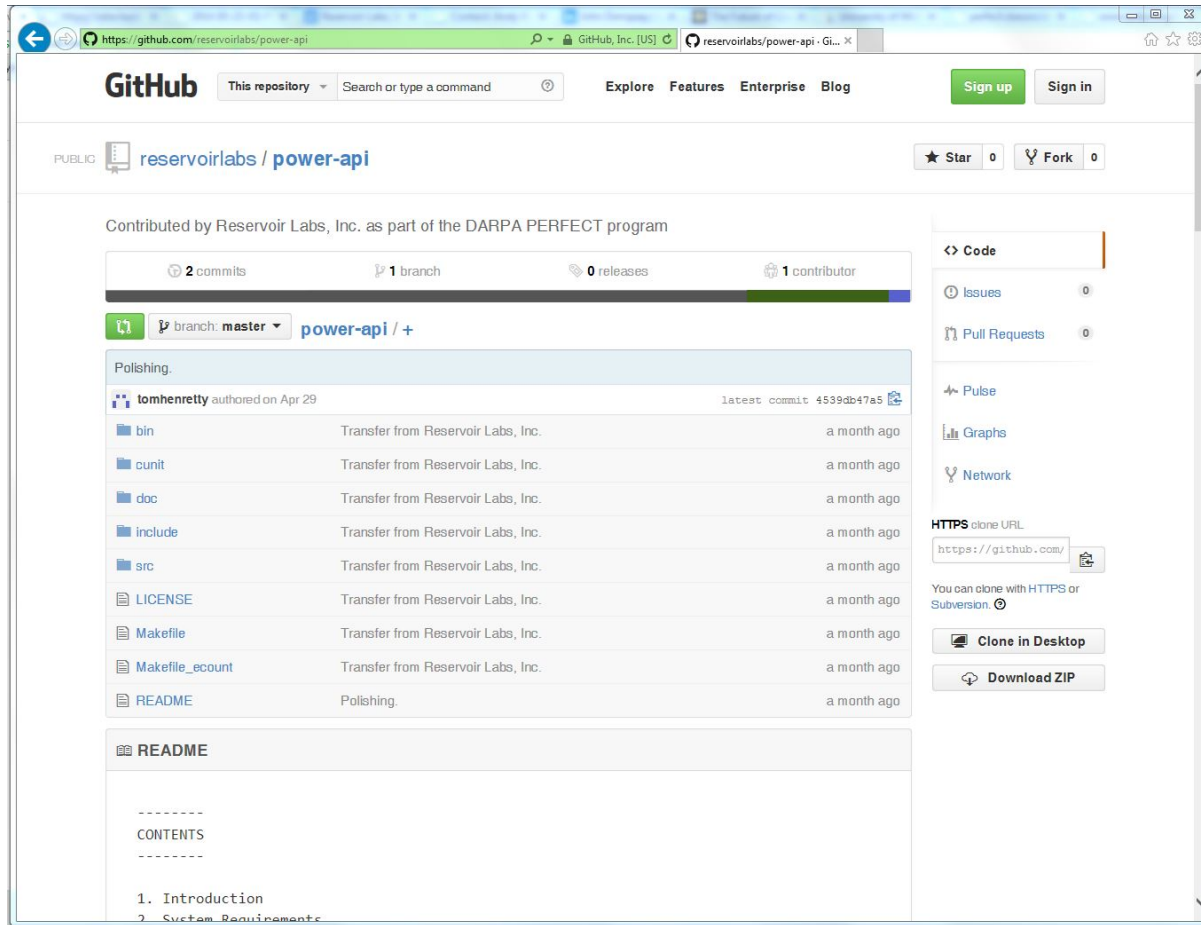
R-Stream performs EPS strategy and statically inserts light-weight runtime power API calls



Betatree representation in R-Stream

```
for(wt=0;wt<=6;wt++){
 for(yt=0;yt<=7;yt++){
  for(xt=0;xt<=7;xt++){
   R-Stream Power API call();
   for(wi=…){
    for(yi=0…){
     for(xi=0…){
      … /* load array out */ }}}
   R-Stream Power API call();
   for(yi=…){
    for(xi=0…){
     for(c=0…){
      … /* load array in */ }}}
   R-Stream Power API call();
   for(wi=0;wi<=7;wi++){
    for(yi=0;yi<=min(…,127);yi++){
     for(xi=0;xi<=min(…,127);xi++){
      for(c=0;c<=31;c++){
       … /* load array filt */
       out_l[8*wt+wi][…][…] += …
}}}}
   R-Stream Power API call();
   for(wi=…){
    for(yi=0…){
     for(xi=0…){
      … /* store array out */ }}}
```
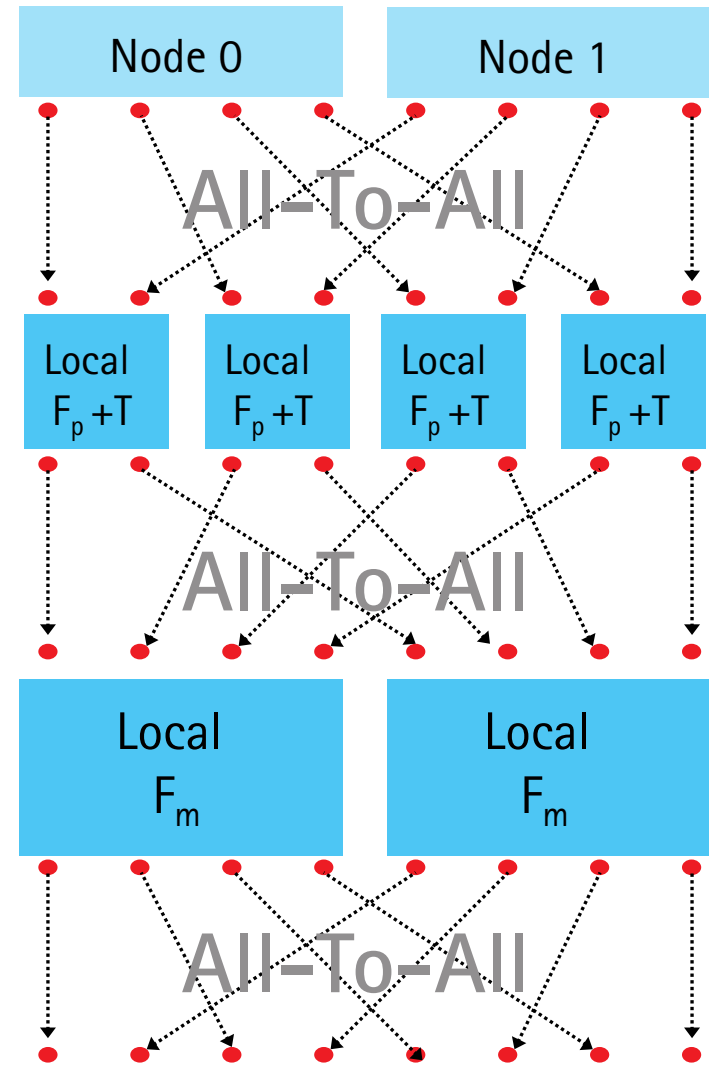
# Power API



https://github.com/reservoirlabs/power-api

# Algorithm Facets

# Traditional Parallel FFT Approach in 1D

Assume N = M*P Points and
Apply The Operator:

- $F_n = (I_p \otimes F_m)(F_p \otimes I_m)\Pi$
- Perform Global Bit Reversal
- Perform Local FFTs and Global Transpose
- Apply Twiddle Factors
- Perform Local FFTs and Global Bit Reversal

- REQUIRES 3 GLOBAL ALL-TO-ALL CALLS!
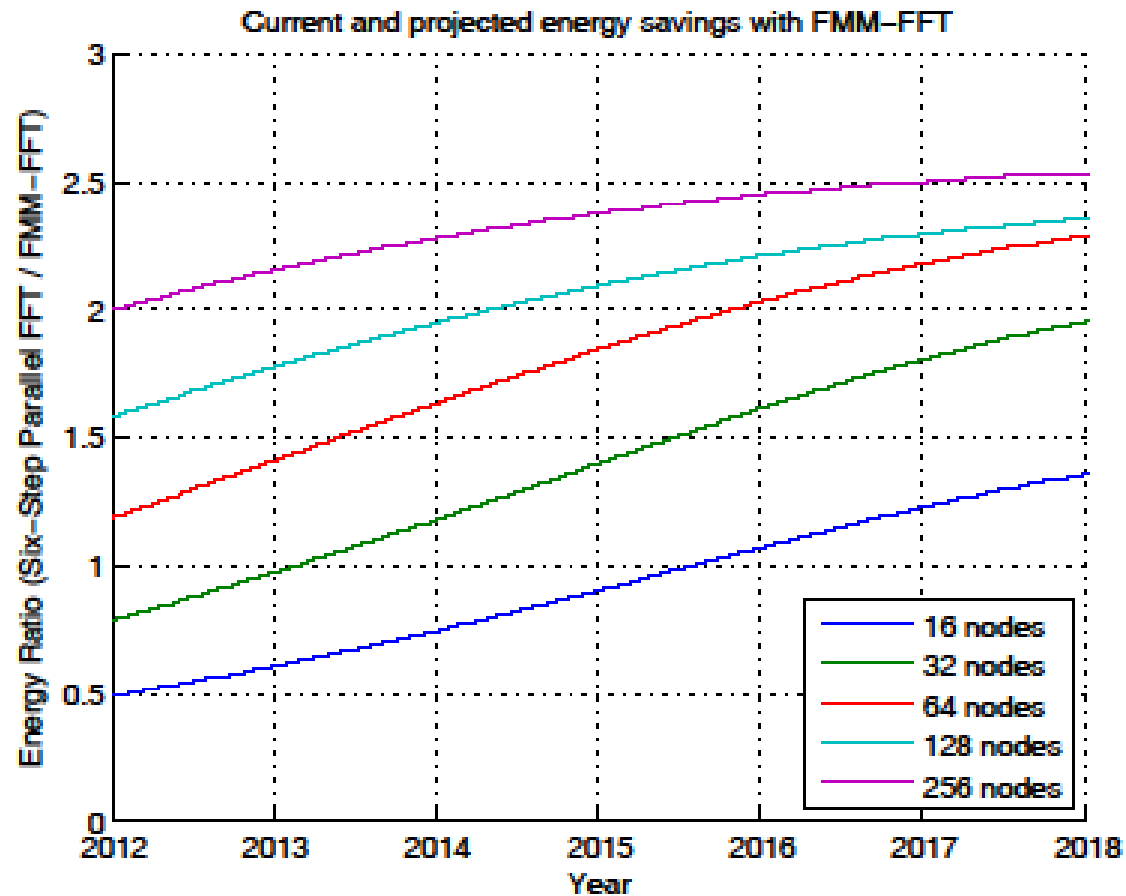
# Alternative Low-Communication FFTs Approximate Computing

Reducing Global Communication from 3 to ~1

- Tang et. al (SC'12) use an oversampling approach while mentioning older approach:

- Edelman et. al. (SIAM J.SciComp'97) has largely been ignored due to previous lack of interest and reliance on FMMs:

- Refactor the operator as $F_n = (I_p \otimes F_m)(F_p \otimes I_m)M\Pi$ with factor matrices $M = diag(I_m, C^1, ..., C^{p-1})$

- $C$ matrices applied with FMM to reduce global communication while arithmetic cost increases. Results (N=$10^7$ points on Intel(R) Xeon(R) CPU X5650 @ 2.67GHz CPUs and QDR InfiniBand):
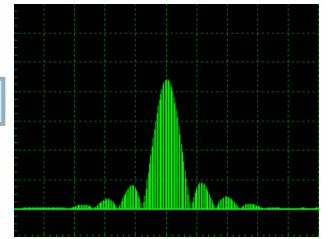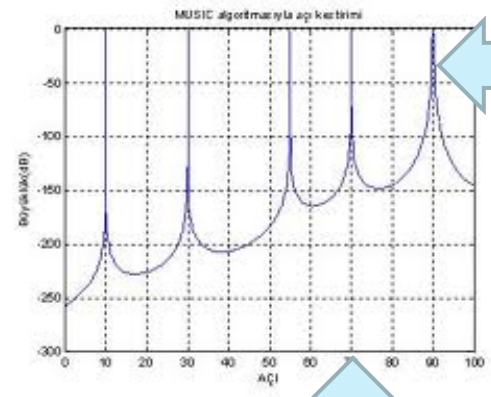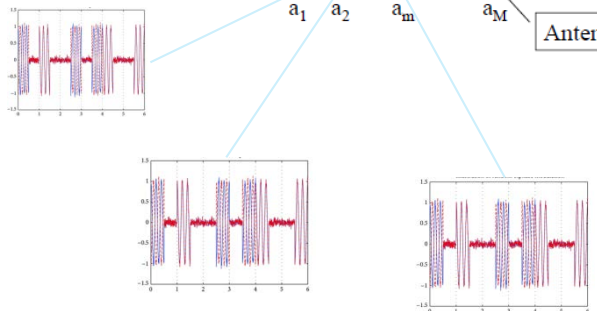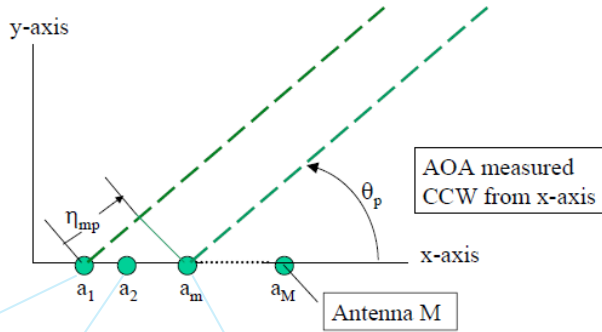
| $NP$ | $T_{FMM}^{ops}$ | $T_{FFT_{FMM}}^{ops}$ | $C_{FMM}^{msgs}$ | $C_{A2A}^{msgs}$ | $C_{TOT}^{msgs}$ | $W_{tot}^{FFT_P}(s)$ | $W_{tot}^{FFT_{FMM}}(s)$ |
|------|------|------|------|------|------|------|------|
| 2 | $4.92e+09$ | $5.92e+09$ | $8.58e+02$ | $4.19e+06$ | $4.20e+06$ | $1.44e+01$ | $3.50e+00$ |
| 4 | $3.33e+09$ | $3.83e+09$ | $2.39e+03$ | $3.15e+06$ | $3.15e+06$ | $3.83e+00$ | $2.58e+00$ |
| 16 | $9.94e+08$ | $1.12e+09$ | $1.05e+04$ | $9.83e+05$ | $9.94e+05$ | $5.74e-01$ | $9.54e-01$ |
| 32 | $5.11e+08$ | $5.74e+08$ | $2.02e+04$ | $5.08e+05$ | $5.28e+05$ | $2.23e-01$ | $4.86e-01$ |
| 64 | $2.60e+08$ | $2.91e+08$ | $3.81e+04$ | $2.58e+05$ | $2.96e+05$ | $1.45e-01$ | $2.96e-01$ |
| 128 | $1.33e+08$ | $1.48e+08$ | $7.07e+04$ | $1.30e+05$ | $2.01e+05$ | $1.14e-01$ | $2.48e-01$ |

- Despite ~10x increase in arithmetic, ~3x reduction in communication results in only 2x slower runtimes

# Energy Savings for FMM-FFT



Current and projected energy savings with FMM-FFT

# Angle of Arrival/Spectrum Sensing Problem



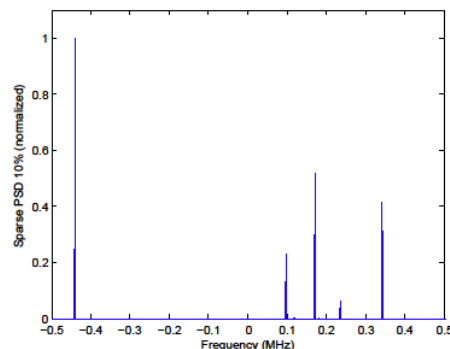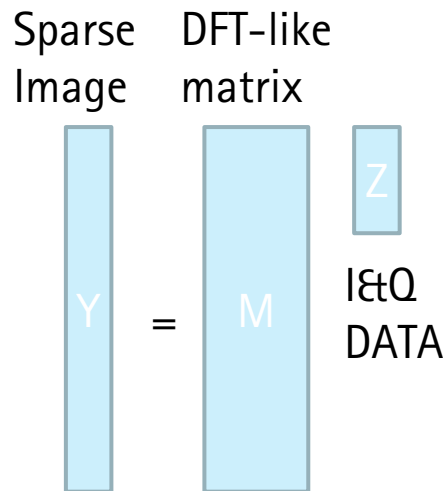Determine spectrum
of each emitter

Locate emitters
in Angle

| Radar Band | Frequency (GHz) | Wavelength (cm) |
|---|---|---|
| Millimeter | 40 to 100 | 0.75 to 0.30 |
| Ka | 26.5 to 40 | 1.1 to 0.75 |
| K | 18 to 26.5 | 1.7 to 1.1 |
| Ku | 12.5 to 18 | 2.4 to 1.7 |
| X | 8 to 12.5 | 3.75 to 2.4 |
| C | 4 to 8 | 7.5 to 3.75 |
| S | 2 to 4 | 15 to 7.5 |
| L | 1 to 2 | 30 to 15 |
| UHF | 0.3 to 1 | 100 to 30 |

Fastest ADC can only sample a fraction of radar band

# Compressive Sensing -  kW to W

Processor 12 kW

Sensor 12 kW



24 kW

## Classical Baseline

Sparse Image    DFT-like matrix

$$Y = M \quad Z$$

I&Q DATA

Sensor 100 W
Processor 100 W



200 W

## Compressive Sensing

# Benchmarking and Evaluation Implications

Hardware programming complexity increasing to save power; suites coded to the metal not portable

Many aspects of new programming complexity can be addressed by automatic tools – consider new tools in suites

- Provides some performance portability
- Input code for benchmarks should be "textbook form"

Algorithm landscape is changing too – expand suites

- Communication-avoiding algorithms
- Sparsity-exploiting algorithms (compressed sensing, etc.)
- Randomized algorithms
- Spectral support solvers

Need to consider system power