# MARK HOROWITZ

## STANFORD UNIVERSITY

# AHA! – VISUAL COMPUTING

# SCALING PROVIDED A GREAT RIDE

**For a 2x scaling**
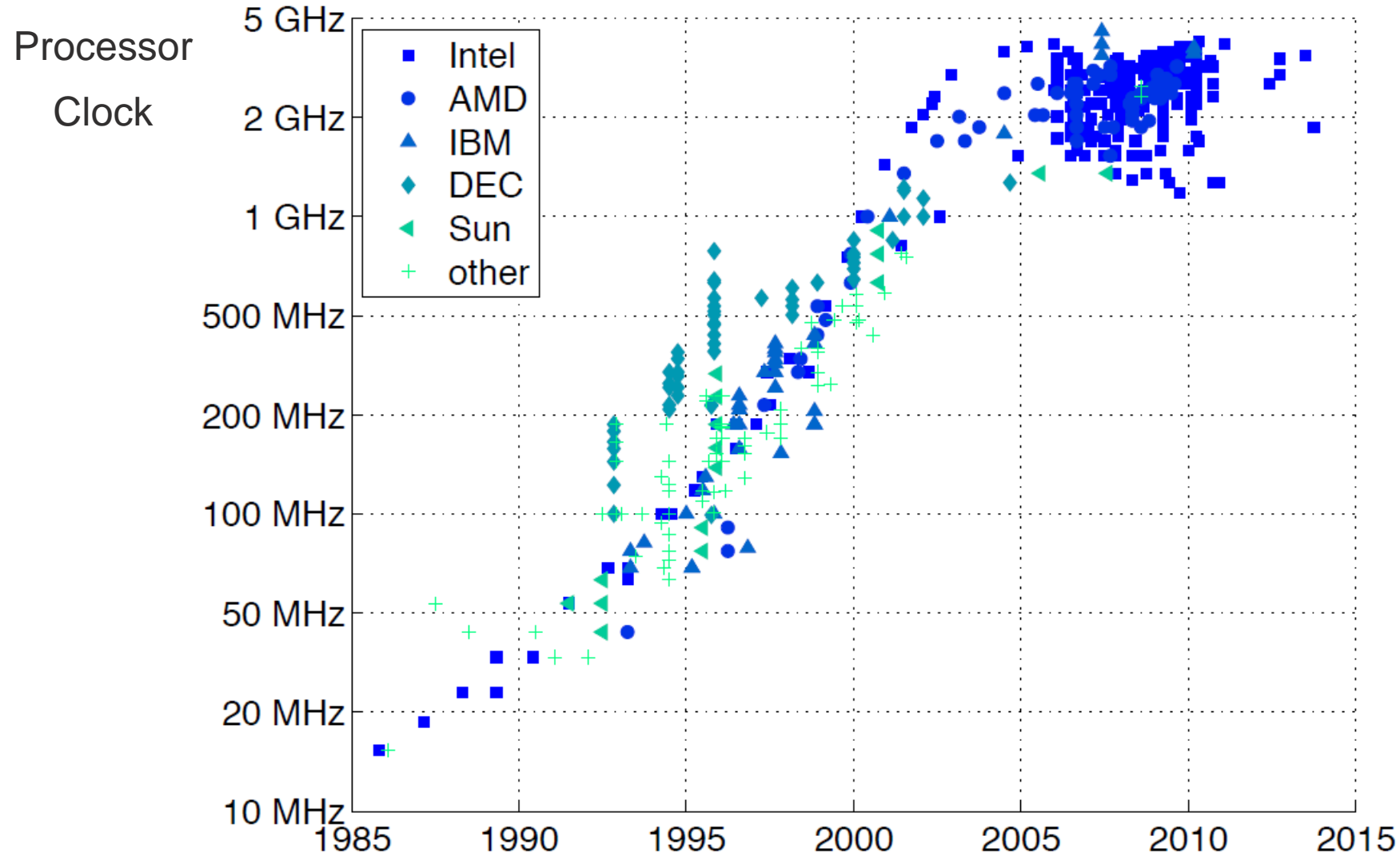
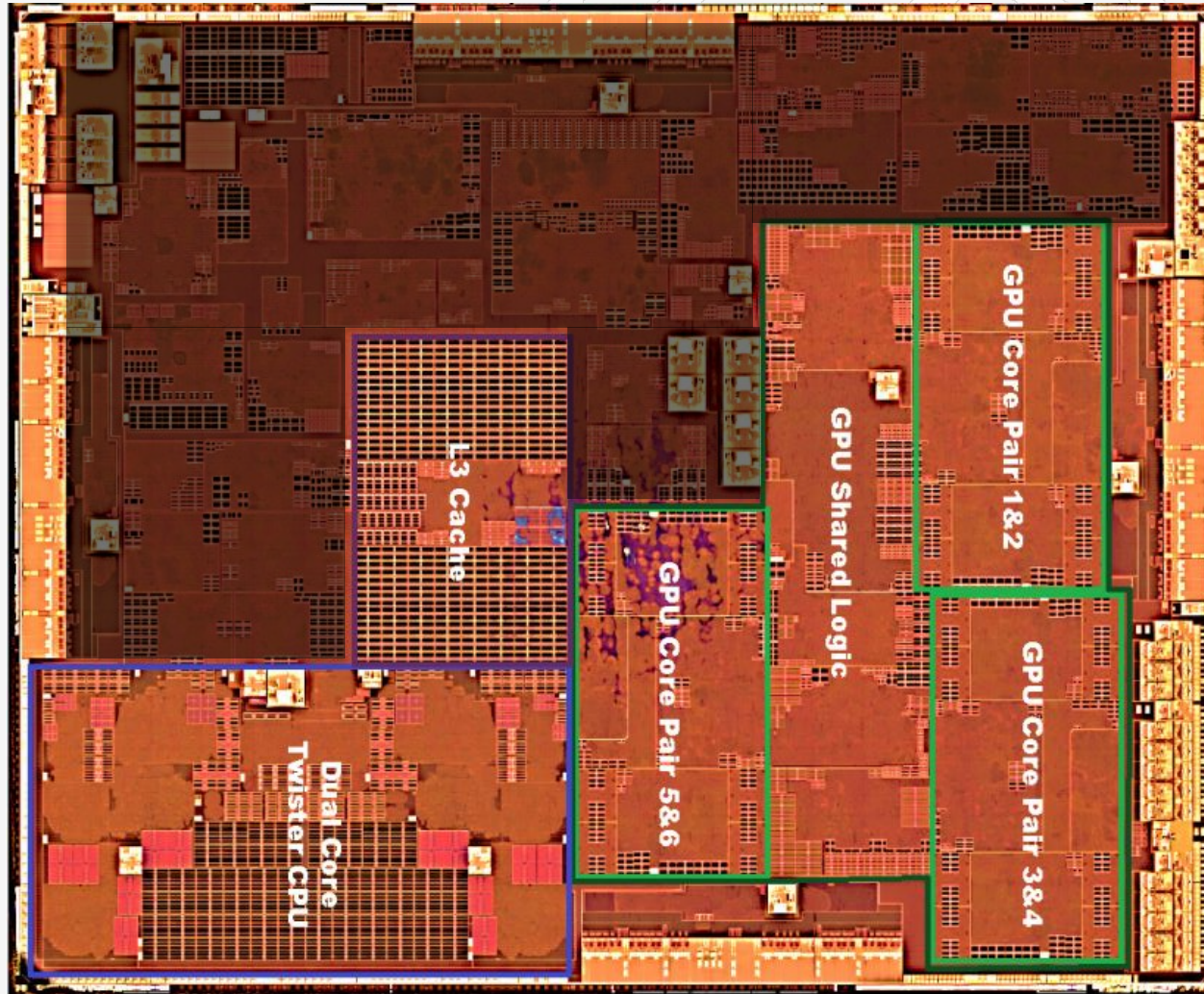**Get 4x more gates,**

**Gates get 2x faster,**

**Energy decrease 8x**

Dennard, JSSC,
pp. 256-268, Oct. 1974

No Exponential is Forever...but We Can Delay 'Forever', Moore ISSCC 2002

# HOUSTON, WE HAVE A PROBLEM

http://cpudb.stanford.edu/

# TO CONTINUE TO SCALE PERFORMANCE



Apple's A9 (2015)

# HOW TO CREATE THESE ACCELERATORS?

- **Study application**

# HOW TO CREATE THESE ACCELERATORS?

- Study application

- **Design hardware**

```
49  //================================================================//
50  // sram-memory interface
51  //================================================================//
52  always_comb begin
53      if (config_en && config_wr) begin
54          // Configuration assumes that 2 * CONFIG_DATA_WIDTH >= BANK_DATA_WIDTH
55          if (CONFIG_DATA_WIDTH * 2 < BANK_DATA_WIDTH)
56              $error("Configuration data width must be at least hal
57          if (config_addr[ADDR_OFFSET-1] == 0) begin
58              // configuring LSB bits
59              sram_to_mem_wen = 1;
60              sram_to_mem_ren = 0;
61              sram_to_mem_cen = 1;
62              sram_to_mem_addr = config_addr[ADDR_OFFSET +: BANK_AD
63              sram_to_mem_data = {{{BANK_DATA_WIDTH-CONFIG_DATA_WID
64              sram_to_mem_bit_sel = {{{BANK_DATA_WIDTH-CONFIG_DATA_
65              config_rd_data = 0;
66          end
67          else begin
68              // configuring MSB bits
69              sram_to_mem_wen = 1;
70              sram_to_mem_ren = 0;
71              sram_to_mem_cen = 1;
72              sram_to_mem_addr = config_addr[ADDR_OFFSET +: BANK_AD
73              sram_to_mem_data = {config_wr_data[BANK_DATA_WIDTH-CO
74              sram_to_mem_bit_sel = {{{BANK_DATA_WIDTH-CONFIG_DATA_
75              config_rd_data = 0;
76          end
77      end
78      else if (config_en && config_rd) begin
79          sram_to_mem_wen = 0;
80          sram_to_mem_ren = 1;
81          sram_to_mem_cen = 1;
82          sram_to_mem_addr = config_addr[ADDR_OFFSET +: BANK_ADDR_W
83          sram_to_mem_data = 0;
84          sram_to_mem_bit_sel = 0;
85          if (config_addr[ADDR_OFFSET-1] == 0) begin
86              config_rd_data = data_out[0 +: CONFIG_DATA_WIDTH];
87          end
88          else begin
89              config_rd_data = data_out[BANK_DATA_WIDTH-1 -: CONFIG
90          end
```

```
113 //================================================================//
114 // configuration
115 //================================================================//
116 integer j, k;
117
118 wire [CONFIG_FEATURE_WIDTH-1:0] config_feature_addr;
119 wire [CONFIG_REG_WIDTH-1:0]     config_reg_addr;
120 reg                             config_en_io_ctrl [`$num_io_channels-1`:0];
121 reg                             config_en_io_int;
122
123 assign config_feature_addr = config_addr[0 +: CONFIG_FEATURE_WIDTH];
124 assign config_reg_addr = config_addr[CONFIG_FEATURE_WIDTH +: CONFIG_REG_WIDTH];
125 always_comb begin
126     for(j=0; j<`$num_io_channels`; j=j+1) begin
127         config_en_io_ctrl[j] = config_en && (config_feature_addr == j);
128     end
129     config_en_io_int = config_en && (config_feature_addr == `$num_io_channels`);
130 end
131
132 always_ff @(posedge clk or posedge reset) begin
133     if (reset) begin
134         switch_sel <= 0;
135     end
136     else begin
137         if (config_en_io_int && config_wr) begin
138             case (config_reg_addr)
139                 0: switch_sel <= config_wr_data;
140             endcase
141         end
142     end
143 end
144
145 always_ff @(posedge clk or posedge reset) begin
146     if (reset) begin
147         for(j=0; j<`$num_io_channels`; j=j+1) begin
148             io_ctrl_mode[j] <= 0;
149             io_ctrl_start_addr[j] <= 0;
150             io_ctrl_num_words[j] <= 0;
151             end
```

# HOW TO CREATE THESE ACCELERATORS?

- Study application

- Design hardware

- **Write software**

# NOT SO SECRET DOWNSIDE

# $100M

# Many Years

# NOT SURPRISING



## It is a waterfall model of design!

# SOFTWARE ISN'T BUILT THAT WAY

- Moved away from that style decades ago

- Enables small teams to build amazing apps

# AGILE DESIGN

# Rapidly iterate on end-to-end system
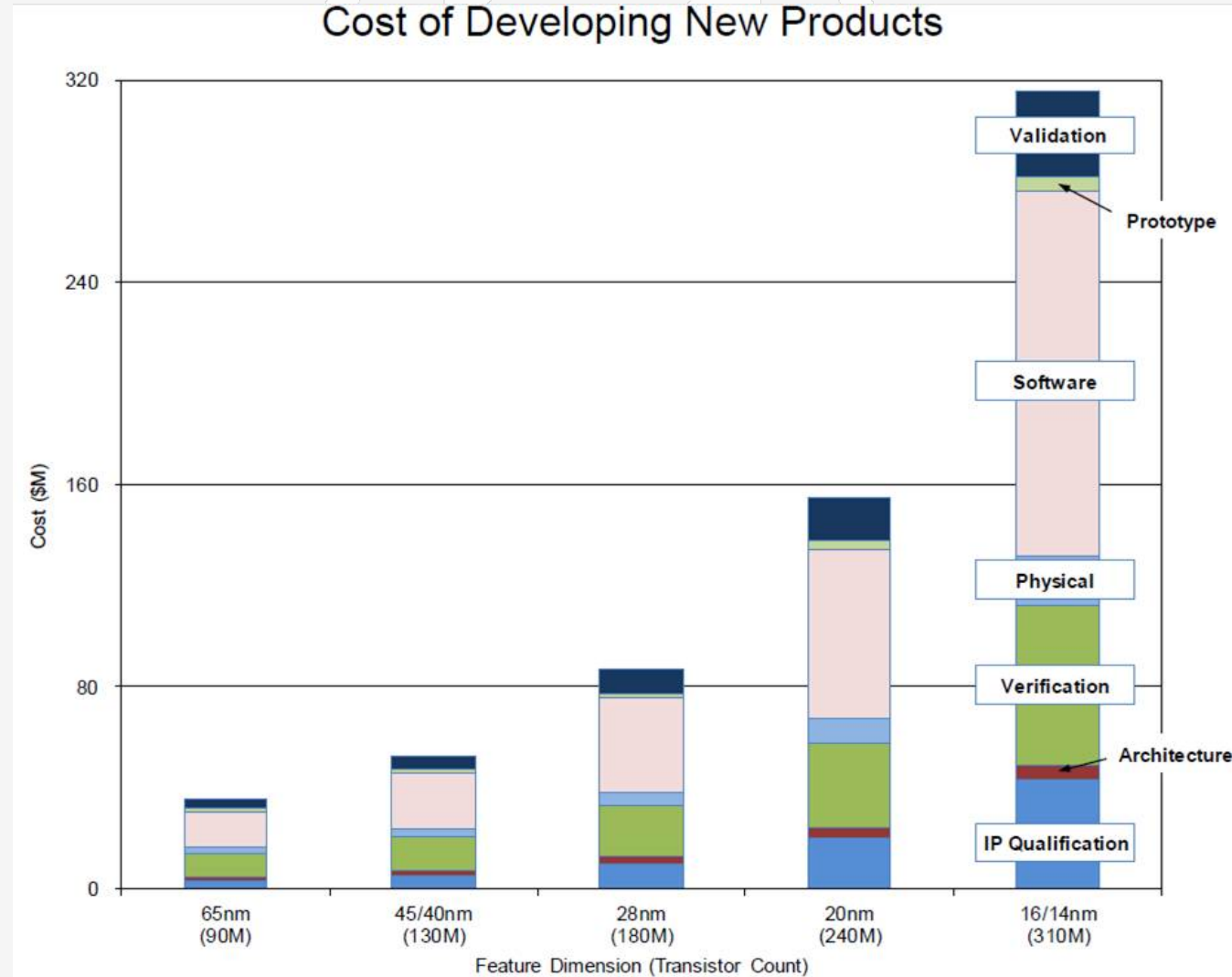
## Learn about real problems, and goals

# AHA!

Agile Hardware

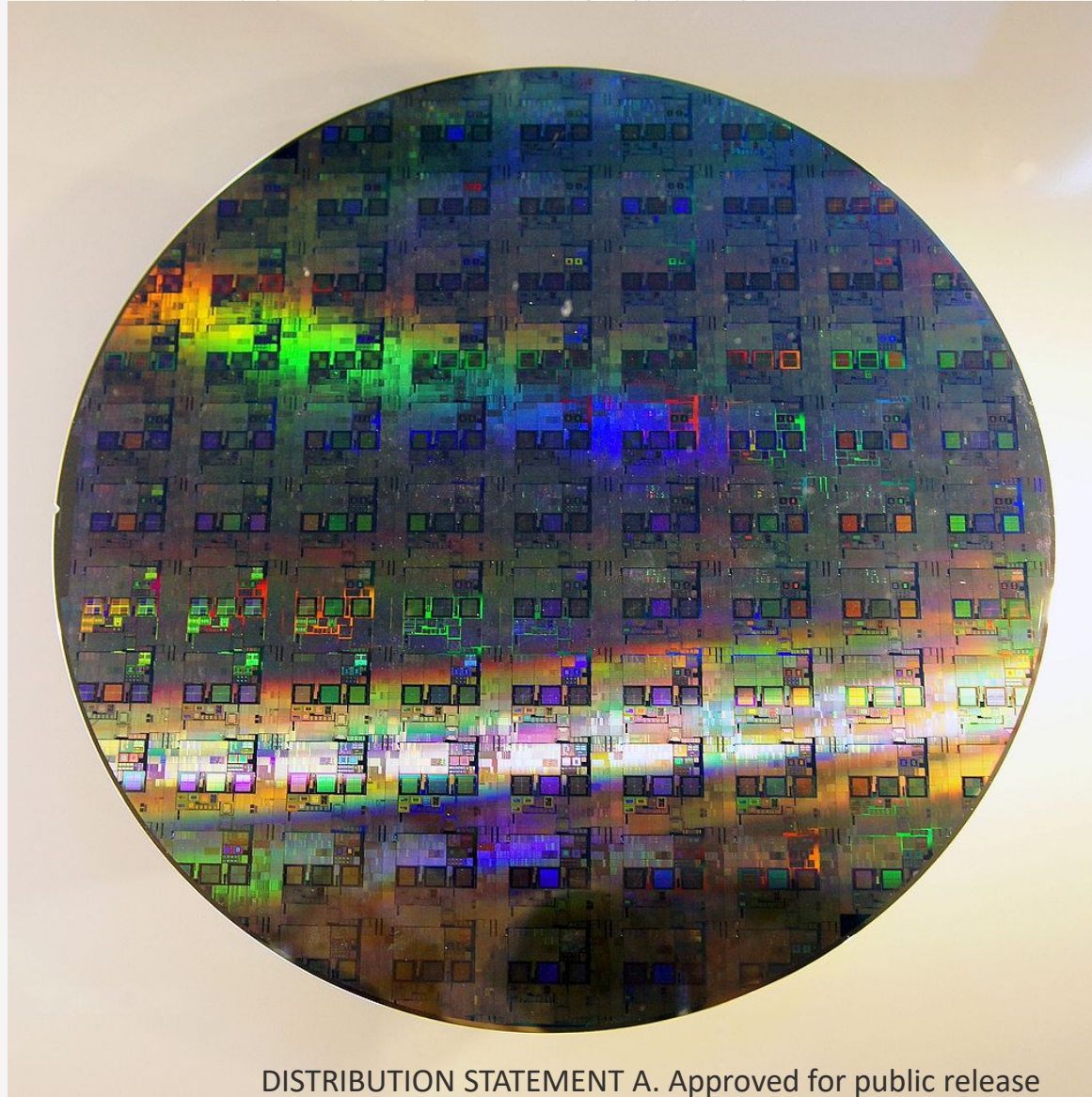# AGILE DESIGN

It is about reuse

It is about clean interfaces

It is about constructors, not instances

It is about productivity

# PRODUCTIVITY IS THE ISSUE IN HARDWARE



Source: IBS

# STILL NEED TO DEAL WITH FABRICATION



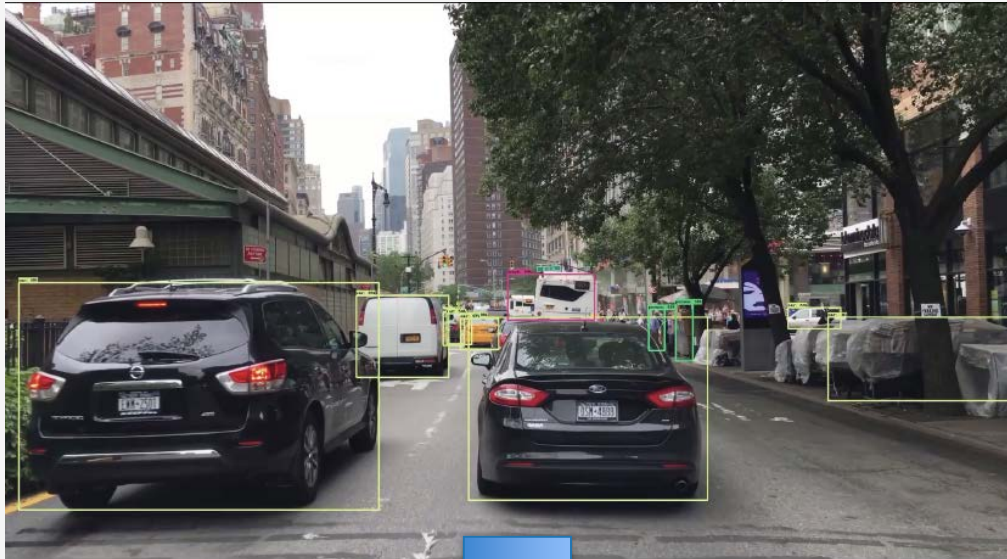DISTRIBUTION STATEMENT A. Approved for public release

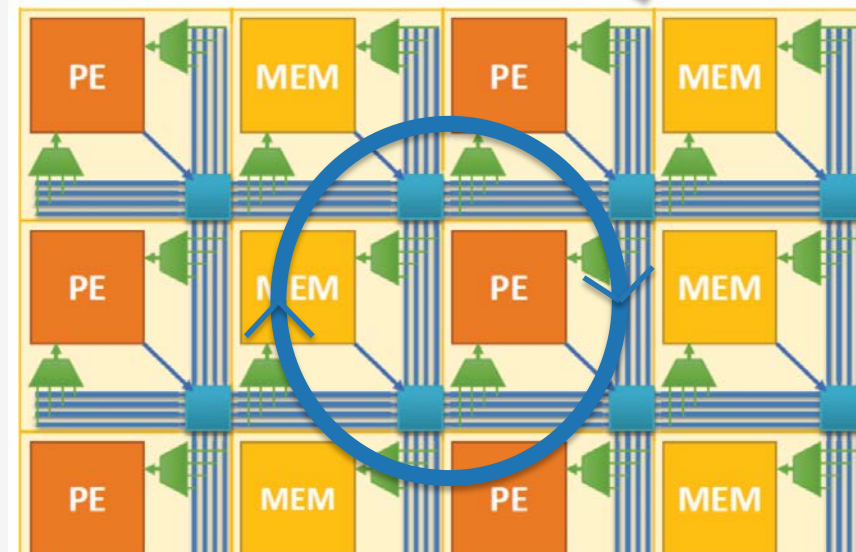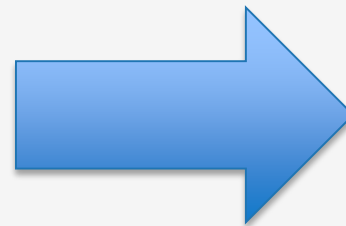# NEED TO EVOLVE THE HARDWARE

- Use a CGRA – a configurable framework
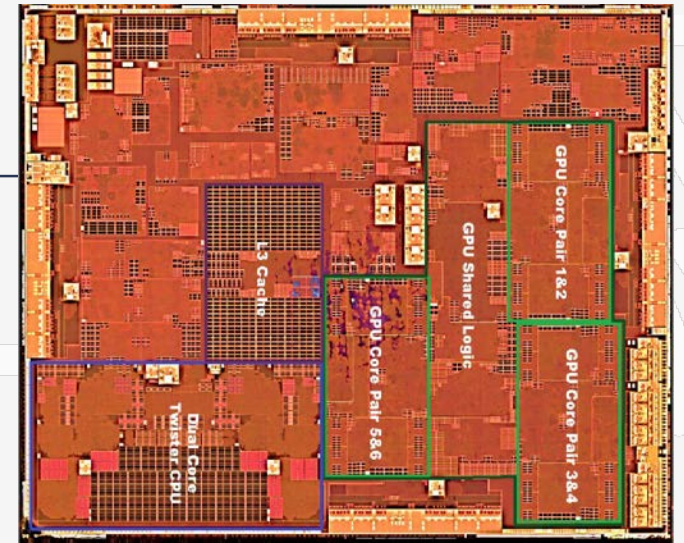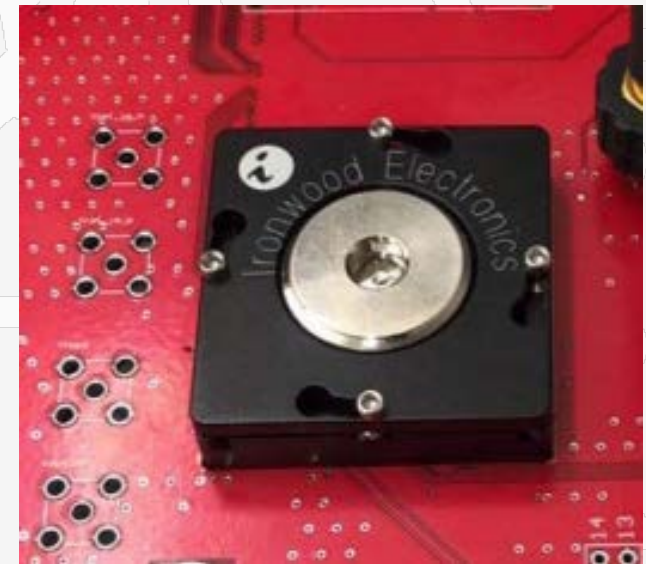
# AHA VISUAL COMPUTING

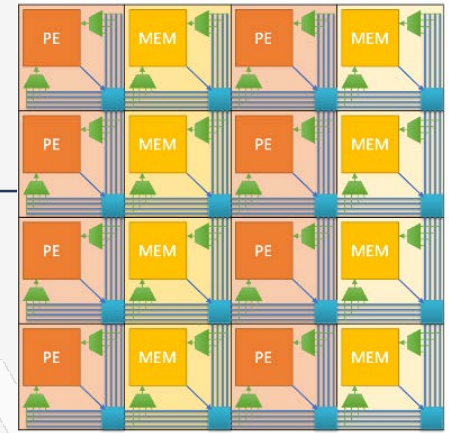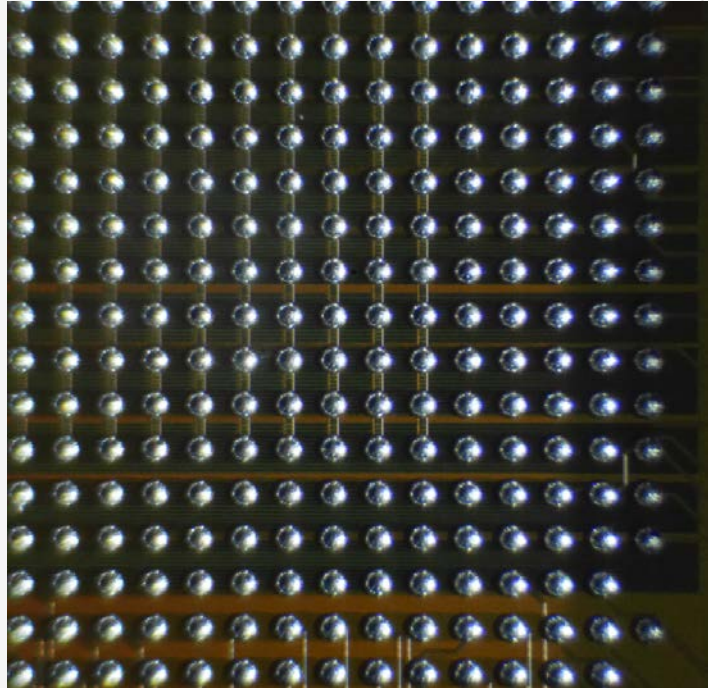**A new way to create DSSoCs**



**Evolves**

**Optimize**
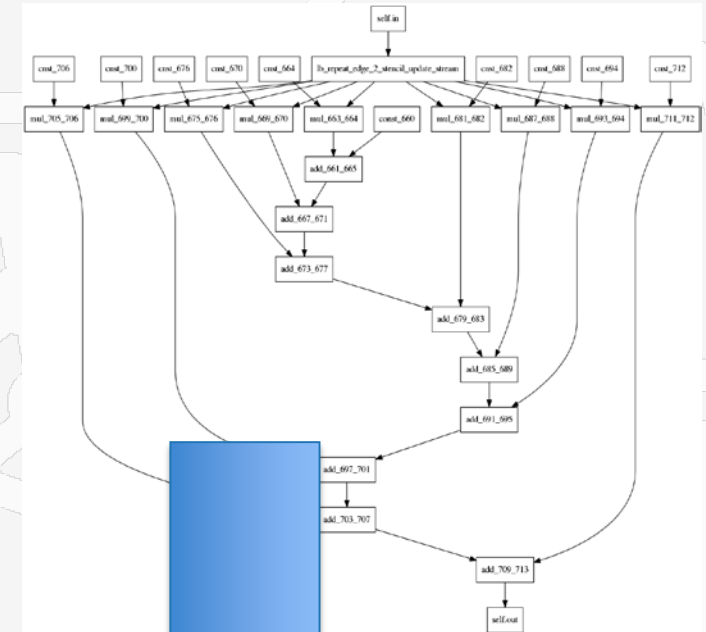
**Compile**

# FIRST GENERATION CGRA

# APPLICATION COMPILATION FLOW
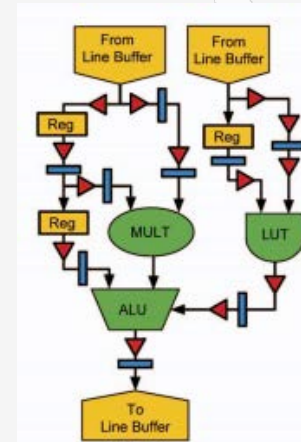
Halide

CoreIR



```
Var x, y, xi, yi, xo, yo; Func conv, out;
RDom win(0,3, 0,3);
kernel(x,y) = {{11,12,13},{14,15,16},{17,18,19}};

// algorithm
conv(x, y) += input(x+win.x, y+win.y) *
                 kernel(win.x, win.y);
out(x, y) = conv(x, y);

// schedule
conv.update(0).unroll(win.x).unroll(win.y);
out.tile(x,y, xo,yo, xi,yi, 62,62).reorder(xi,yi, xo,yo);
conv.linebuffer();
```
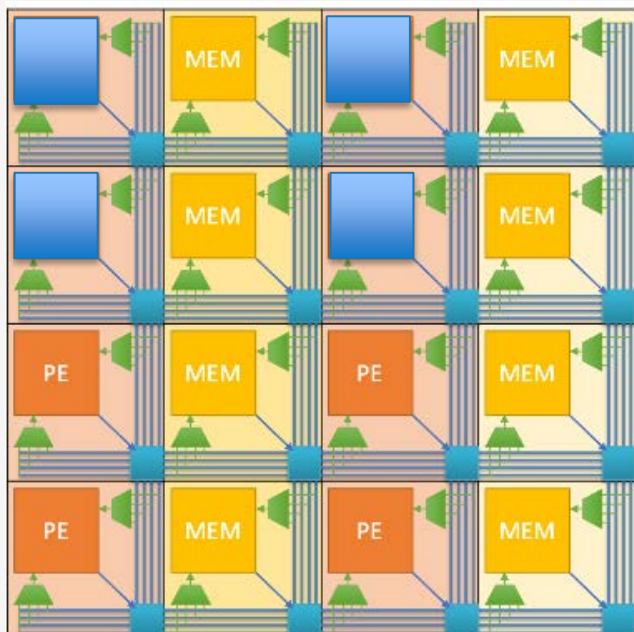
Compiler

Place
& Route
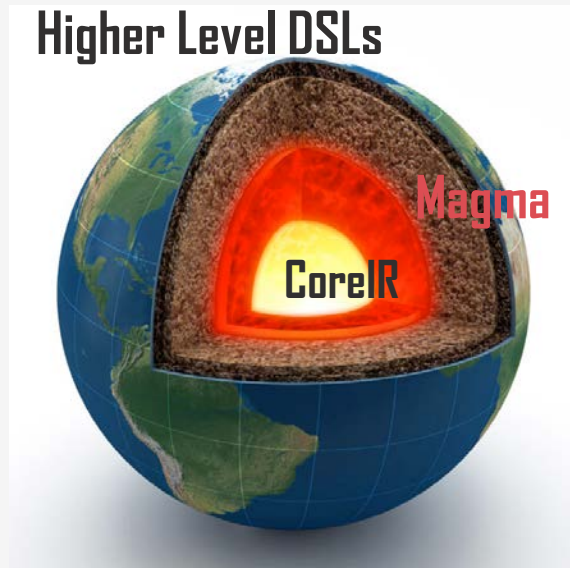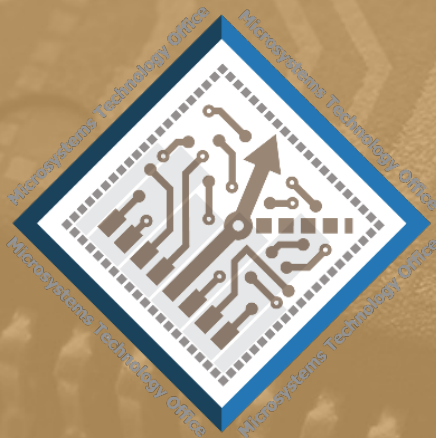
Mapping

# MAINTAINING THE FLOW THROUGH CHANGES

- Many tools need to know about your design
  - You are building a "world"



Image source: Jesse Varner - Raft Lake and the Little Wind River Valley beyond, https://www.flickr.com/photos/jvarner/136566920/in/photostream/

# DSSOC DOESN'T NEED TO BE EXPENSIVE

- One just needs to think about the problem differently

- We have already created one working chip/system using this flow

- And have the next generation of the system working

Stay Tuned for Future Results …