

---

**E225C – Lecture 3**  
**System on a Chip Design**

**Bob Brodersen**

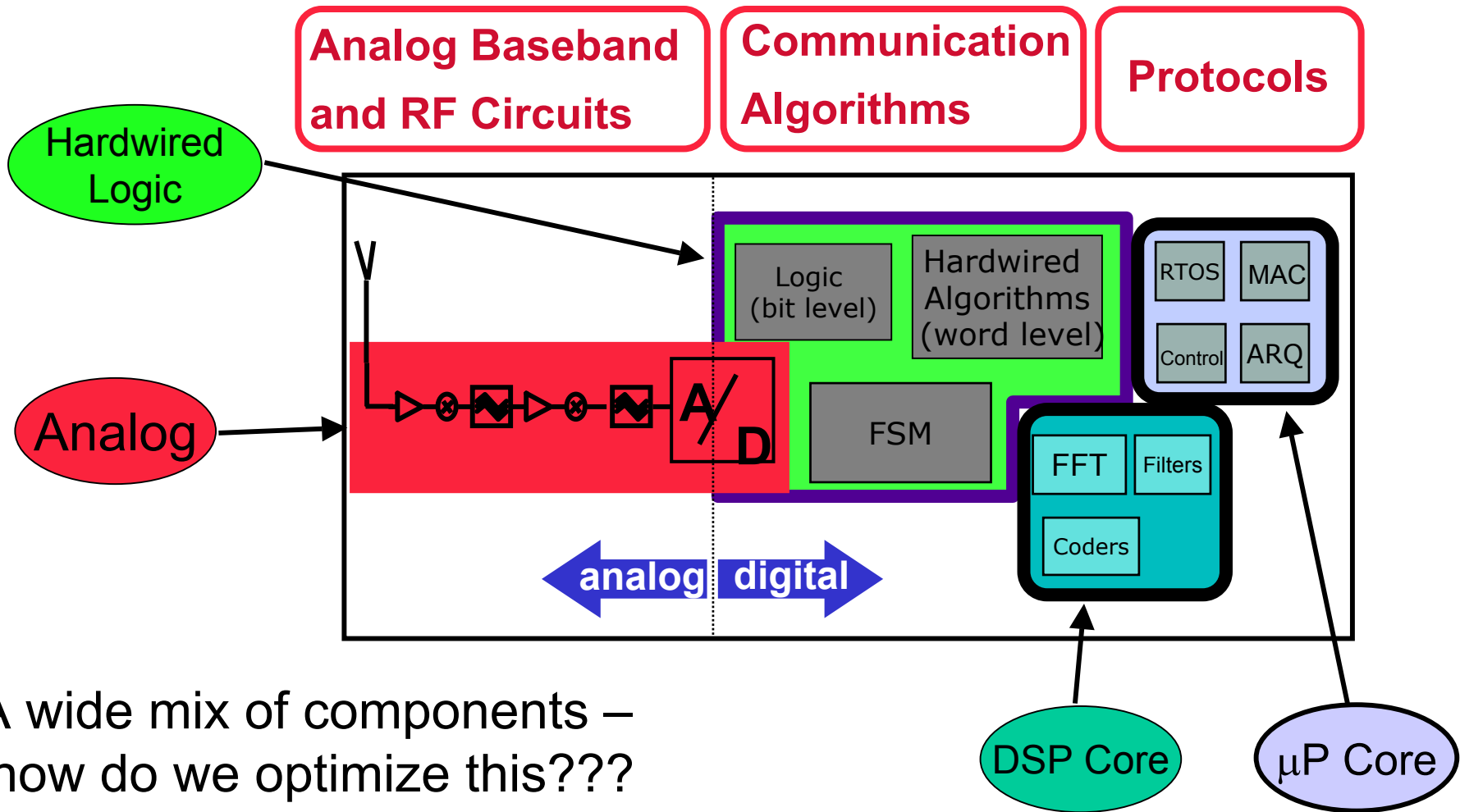
# What is an SoC?

---

Let me define what I think it is....

“A chip designed for “complete” system functionality that incorporates a heterogeneous mix of processing and computation architectures”

# A Wireless System – Typical SOC Design



A wide mix of components –  
how do we optimize this???

# An SOC Design Flow with Prototyping

*Algorithm/flexibility  
evaluation*

Initial System Description  
(Floating point Matlab/Simulink)  
**Determine *Flexibility Requirements***

*Digital delay,  
area and  
energy estimates  
& effect of analog  
impairments*

Architecture/algorithm Description  
with Hardware Constraints (Fixed point Simulink,  
FSM Control in Stateflow)

*Common test vectors,  
and hardware description of  
net list and modules*

Real-time Emulation  
(BEE FPGA Array)

Automated AISC Generation  
(Chip-in-a-Day flow)

# The Issues I am Going to Address

---

- How much flexibility is needed and how best to include it...
- A single system description including interaction between the analog and digital domains
- “Realtime” SOC prototyping
- Automated ASIC design flow

# Flexibility

---

- Determining how much to include and how to do it in the most efficient way possible
- Claims (to be shown)
  - » There are good reasons for flexibility
  - » The “cost” of flexibility is orders of magnitude of inefficiency over an optimized solution
  - » There are many different ways to provide flexibility

# Good reasons for flexibility

---

- One design for a number of SoC customers – more sales volume
- Customers able to provide added value and uniqueness
- Unsure of specification or can't make a decision
- Backwards compatibility with debugged software
- Risk, cost and time of implementing hardwired solutions

*Important to note: these are business, not technical reasons*

# So, what is the cost of flexibility?

---

We need technical metrics that we can look to compare flexible and non-flexible implementations

- A power metric because of thermal limitations
- An energy metric for portable operation
- A cost metric related to the area of the chip
- Performance (computational throughput)

Lets use metrics normalized to the amount of computation being performed – so now lets define computation



# Definitions...

## Computation

- Operation = OP=algorithmically interesting computation (i.e. multiply, add, delay)
- MOPS = Millions of OP's per Second
- $N_{op}$  = Number of parallel OP's in **each** clock cycle

## Power

- $P_{chip} = \text{Total power of chip} = A_{chip} * C_{sw} * (V_{dd})^2 * f_{clk}$
- $C_{sw} = \text{Switched Capacitance/mm}^2$   
 $= P_{chip} / (A_{chip} * V_{dd}^2 * f_{clk})$

## Area

- $A_{chip} = \text{Total area of chip}$
- $A_{op} = \text{Average area of each operation} = A_{chip} / N_{op}$

# Energy Efficiency Metric: MOPS/mW

How much computing (number of operations) can we do with a finite energy source (e.g. battery)?

$$\begin{aligned}\text{Energy Efficiency} &= \frac{\text{Number of useful operations}}{\text{Energy required}} \\ &= \frac{\text{\# of Operations}}{\text{NanoJoule}} = \text{OP/nJ} \\ &= \frac{\text{OP/Sec}}{\text{NanoJoule/Sec}} = \frac{\text{MOPS}}{\text{mW}} \\ &= \text{Power Efficiency}\end{aligned}$$

# Energy and Power Efficiency

---

$$OP/nJ = MOPS/mW$$

Interestingly the energy efficiency metric for energy constrained applications ( $OP/nJ$ ) for a fixed number of operations is the same as that for thermal (power) considerations when maximizing throughput ( $MOPS/mW$ ).

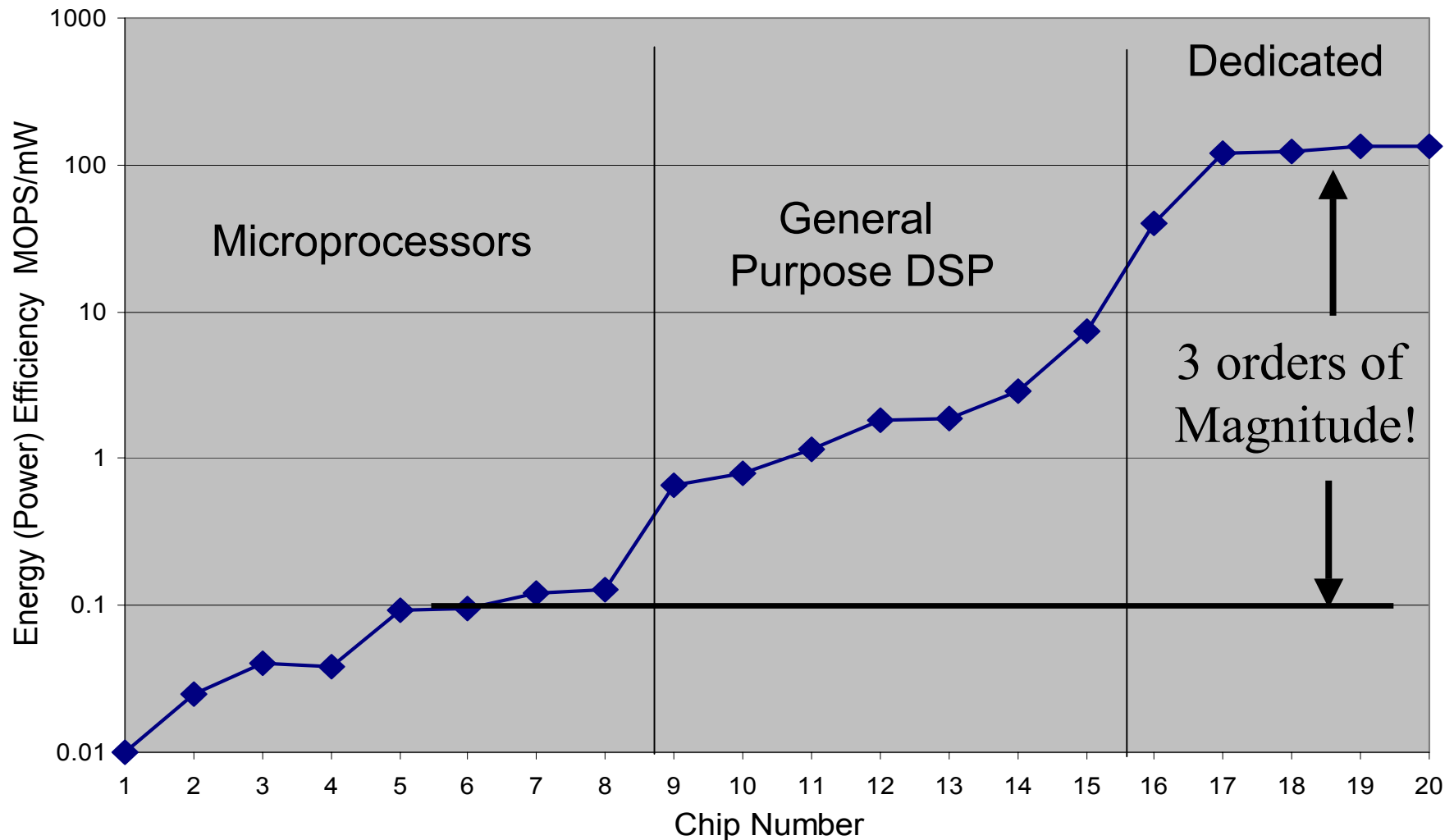
So lets look at a number of chips to see how these efficiency numbers compare

# ISSCC Chips (.18 $\mu$ -.25 $\mu$ )

Chip #	Year	Paper	Description
1	1997	10.3	$\mu$ P - S/390
2	2000	5.2	$\mu$ P - PPC (SOI)
3	1999	5.2	$\mu$ P - G5
4	2000	5.6	$\mu$ P - G6
<b>Microprocessors</b>			
5	2000	5.1	$\mu$ P - Alpha
6	1998	15.4	$\mu$ P - P6
7	1998	18.4	$\mu$ P - Alpha
8	1999	5.6	$\mu$ P - PPC
9	1998	18.6	DSP - StrongArm
10	2000	4.2	DSP - Comm

Chip #	Year	Paper	Description
11	1998	18.1	DSP -Graphics
12	1998	18.2	DSP - Multimedia
<b>DSP's</b>			
13	2000	14.6	DSP - Multimedia
14	2002	22.1	DSP - Mpeg Decoder
15	1998	18.3	DSP - Multimedia
16	2001	21.2	Encryption Processor
17	2000	14.5	Hearing Aid Processor
<b>Dedicated</b>			
18	2000	4.7	FIR for Disk Read Head
19	1998	2.1	MPEG Encoder
20	2002	7.2	802.11a Baseband

# Energy Efficiency (MOPS/mW or OP/nJ)



# What does the low efficiency really mean?

The basic processor architecture puts our circuits at the very limit of failure...

***What happens  
when the  
CPU cooler is  
removed?***



*[www.tomshardware.de](http://www.tomshardware.de)  
[www.tomshardware.com](http://www.tomshardware.com)*

# Why such a big difference?

Lets look at the components of MOPS/mW.

The operations per second:

$$\text{MOPS} = f_{\text{clk}} * N_{\text{op}}$$

The power:

$$P_{\text{chip}} = A_{\text{chip}} * C_{\text{sw}} * (V_{\text{dd}})^2 * f_{\text{clk}}$$

The ratio ( $\text{MOPS}/P_{\text{chip}}$ ) gives the MOPS/mW

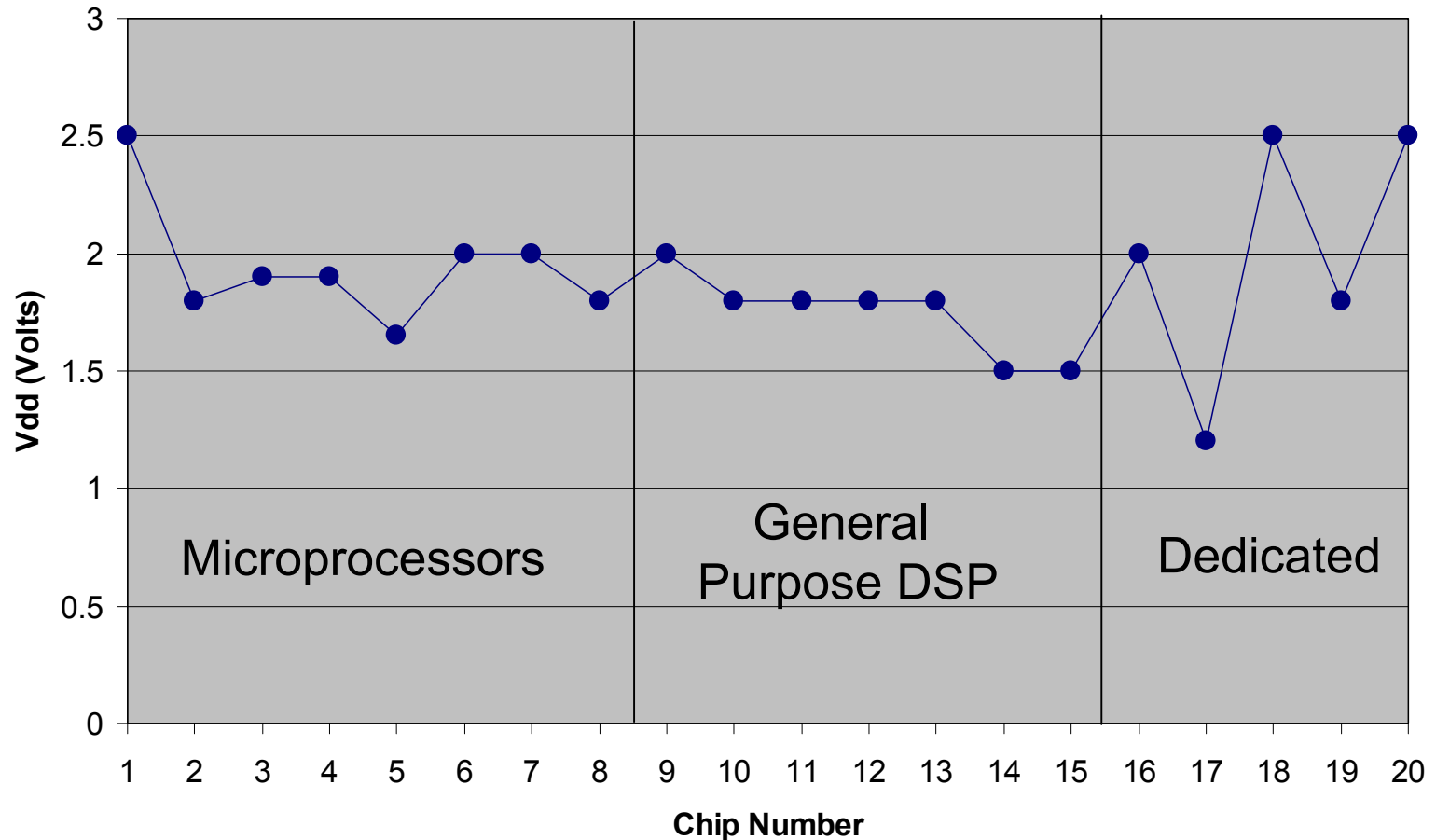
$$= (f_{\text{clk}} * N_{\text{op}}) / A_{\text{chip}} * C_{\text{sw}} * (V_{\text{dd}})^2 * f_{\text{clk}}$$

Simplifying,

$$\text{MOPS/mW} = 1 / (A_{\text{op}} * C_{\text{sw}} * V_{\text{dd}}^2)$$

So lets look at the 3 components –  $V_{\text{dd}}$ ,  $C_{\text{sw}}$  and  $A_{\text{op}}$

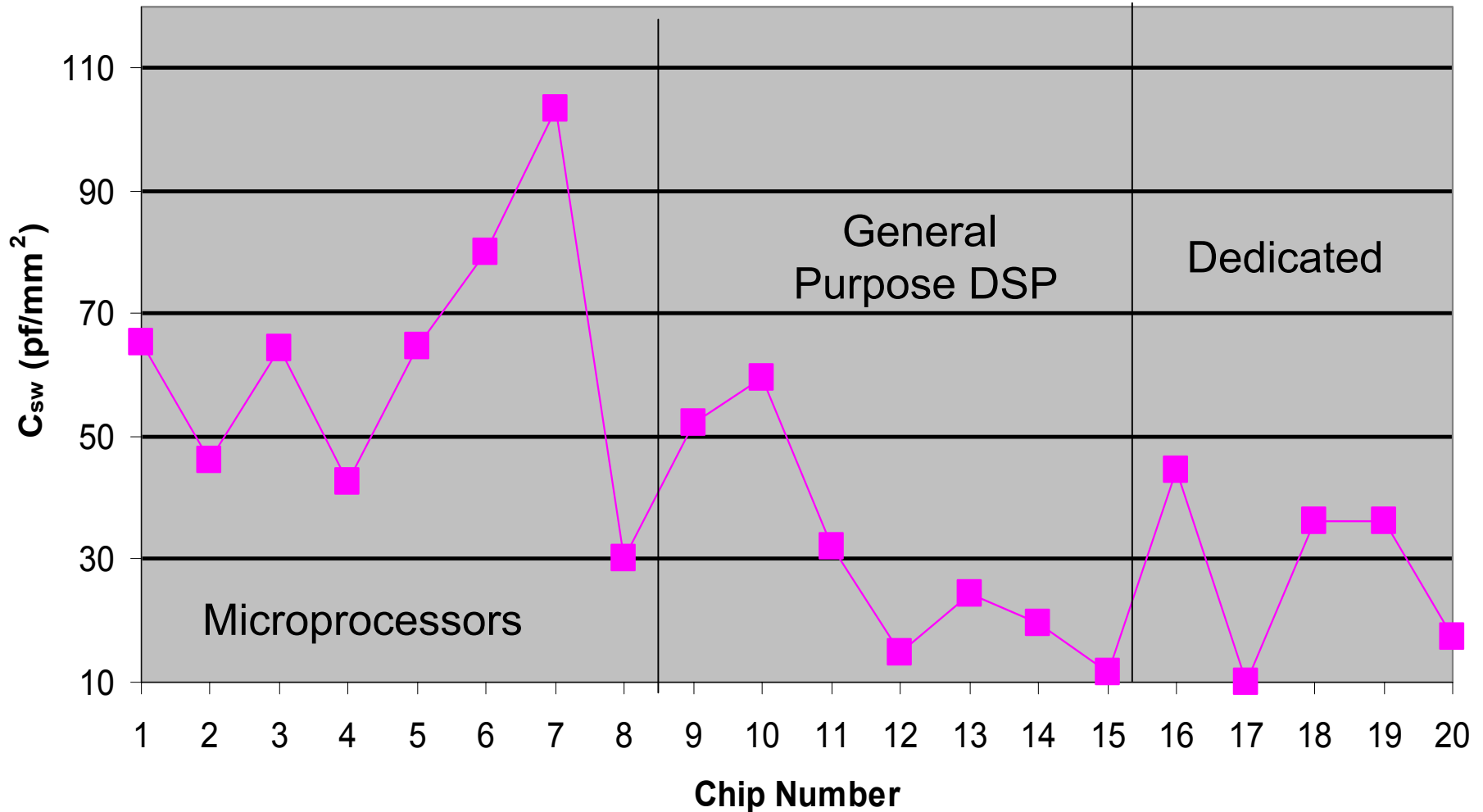
# Supply Voltage, $V_{dd}$



Supply voltage isn't the cause of the difference,  
actually a bit higher for the dedicated chips



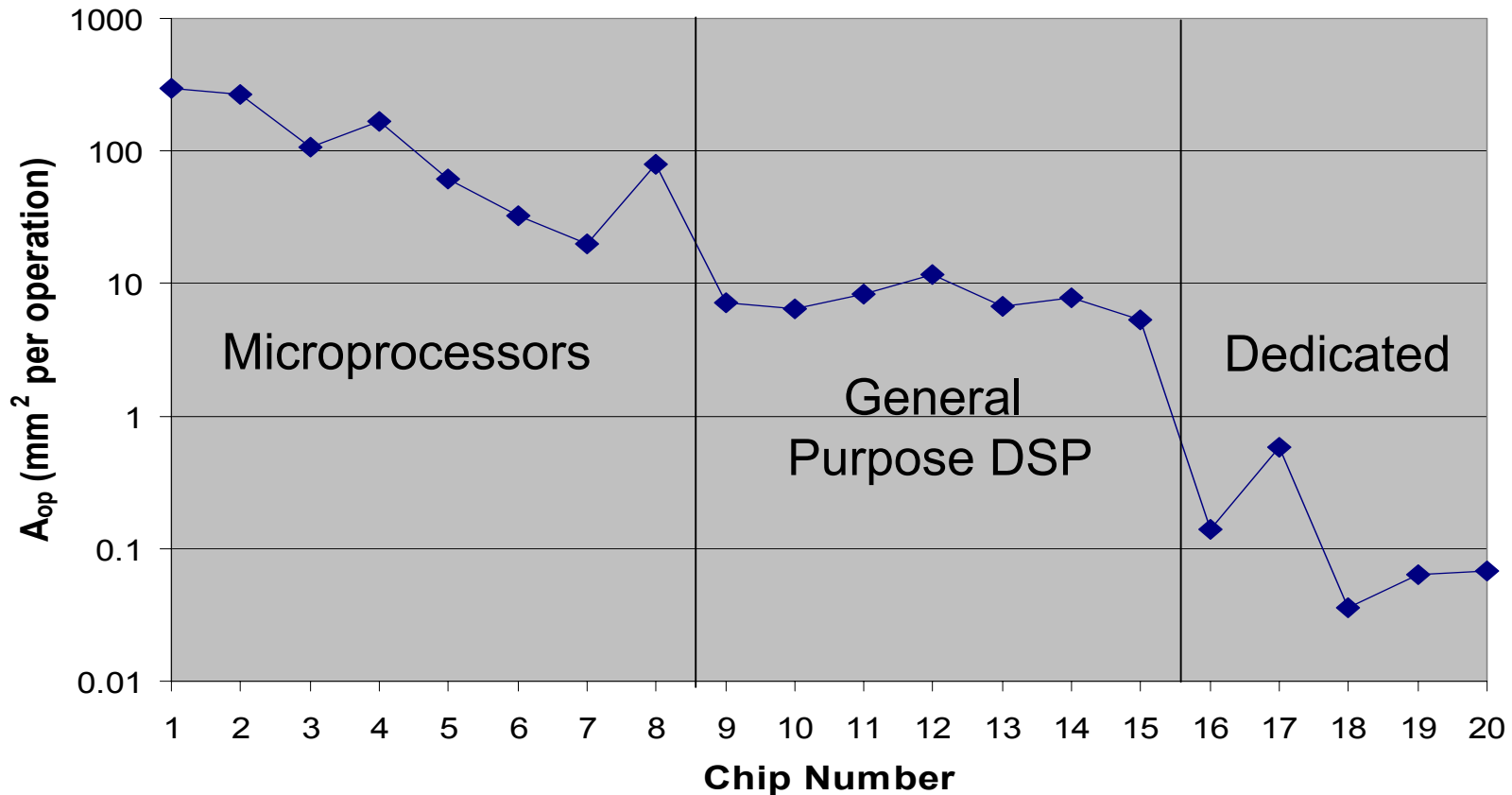
# Switched Capacitance, $C_{sw}$ (pF/mm<sup>2</sup>)



$C_{sw}$  is lower for dedicated, but only by a factor of 2 to 3

# $A_{op} = \text{Area per operation } (A_{chip}/N_{op})$

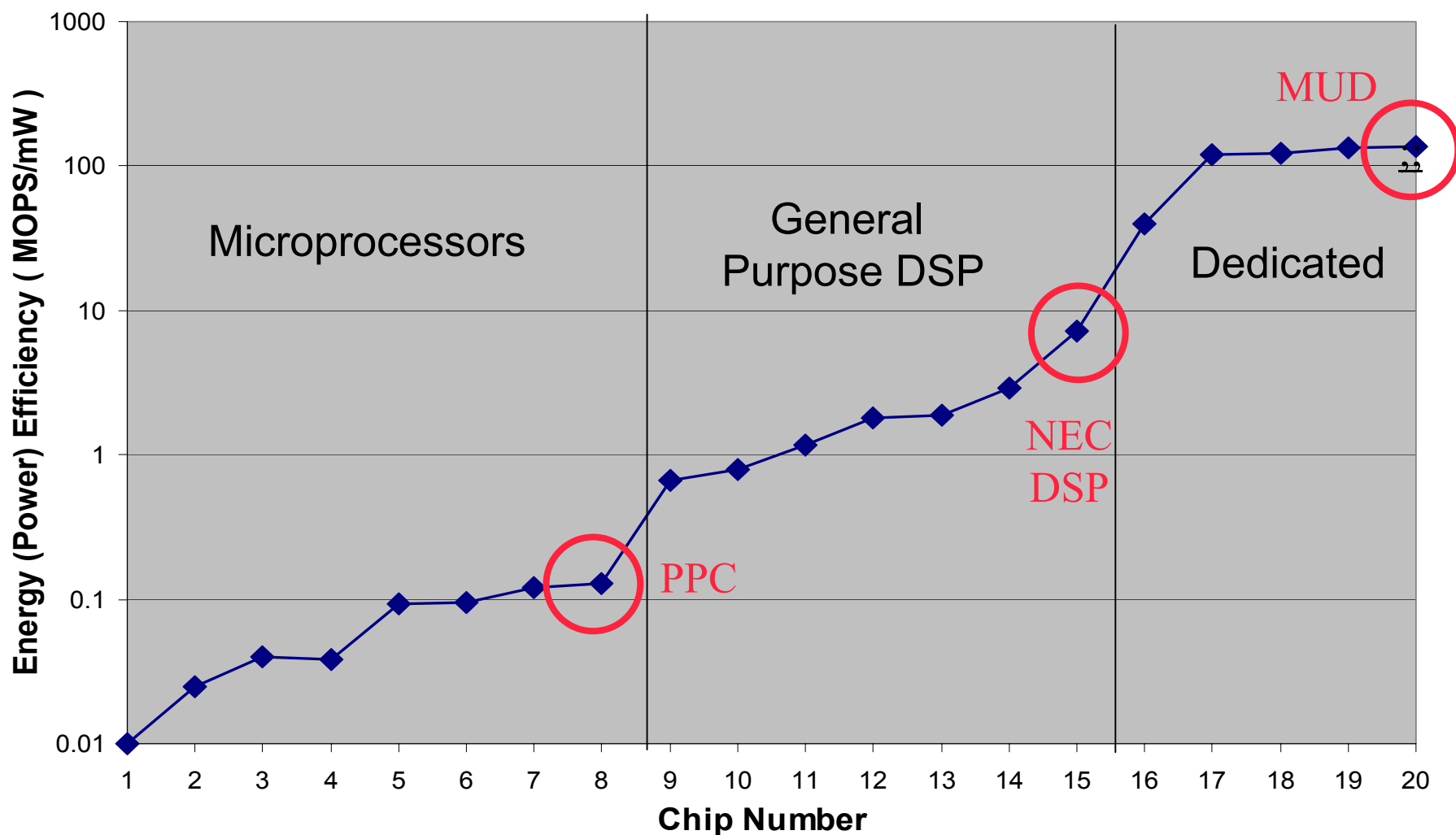
$MOPS/mW = 1/(A_{op} * C_{sw} * V_{dd}^2)$  ;  $A_{op} = A_{chip}/N_{op}$



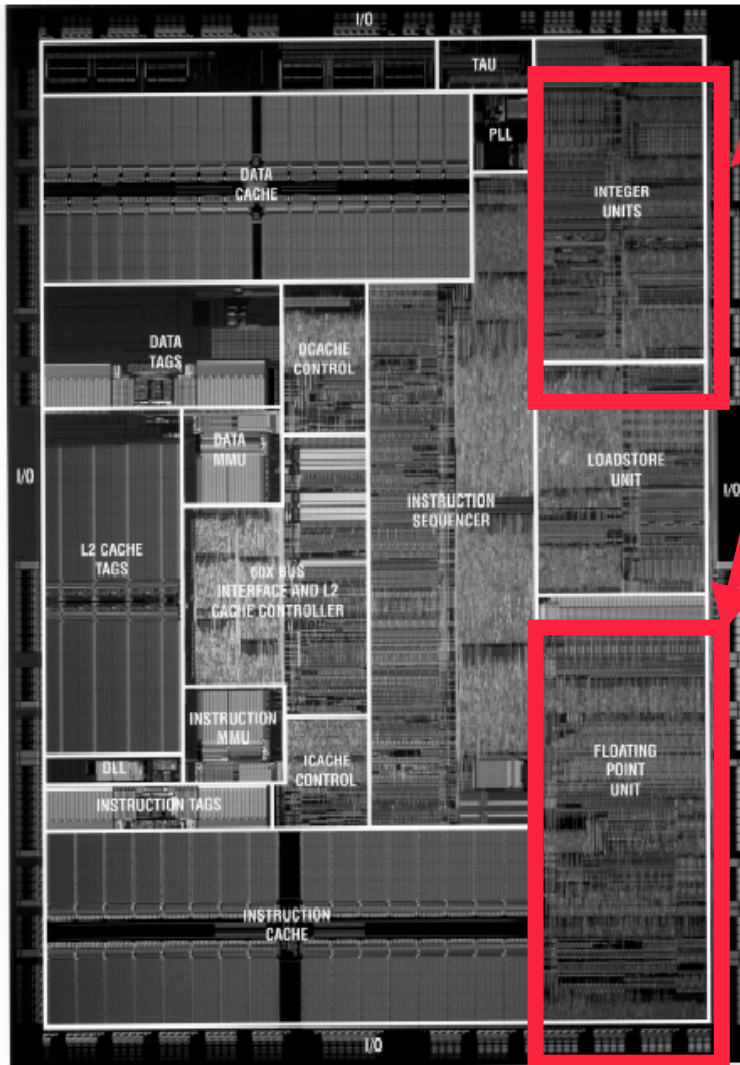
Here is the one that explains the difference, lower due to more parallelism (higher  $N_{op}$ ) in a smaller chip area (less overhead)

# Lets look at some chips to actually see the different architectures

We'll look at one from each category...



# Microprocessor: MOPS/mW=.13



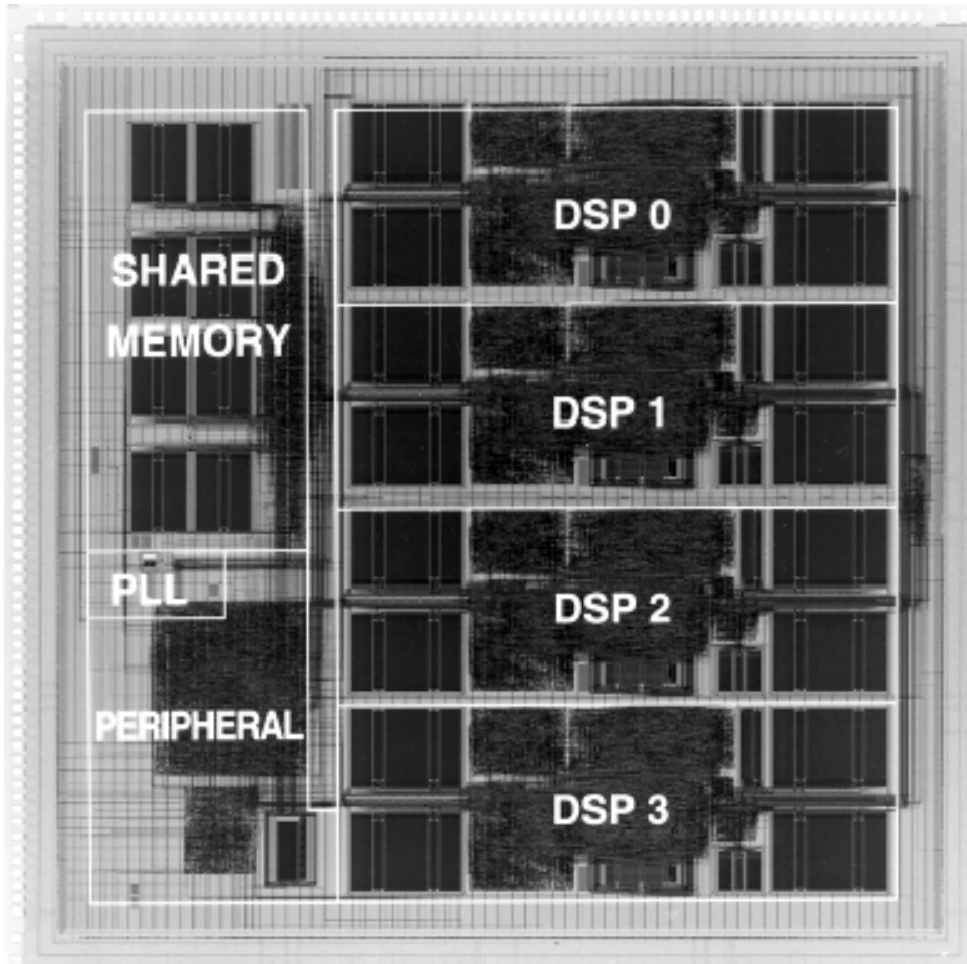
The only circuitry which supports “useful operations”  
All the rest is overhead to support the time multiplexing

$$N_{op} = 2$$
$$f_{clock} = 450 \text{ MHz (2 way)}$$
$$= 900 \text{ MIPS}$$

Two operations each clock cycle, so  
 $A_{op} = A_{chip}/2 = 42 \text{ mm}^2$

Power = 7 Watts

# DSP: MOPS/mW=7



Same granularity (a datapath), more parallelism

4 Parallel processors  
(4 ops each)

$$N_{op} = 16$$

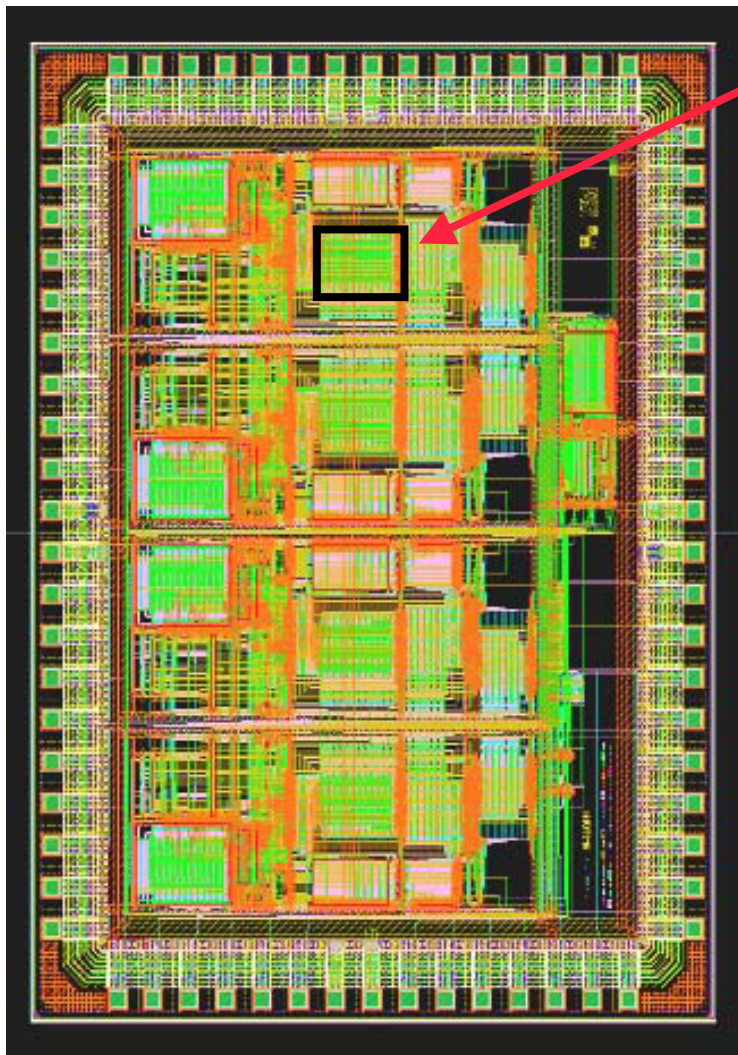
50 MHz clock  
=> 800 MOPS

Sixteen operations  
each clock cycle, so

$$A_{op} = A_{chip}/16 = 5.3\text{mm}^2$$

Power = 110 mW.

# Dedicated Design: MOPS/mW=200



Complex  
mult/add  
(8 ops)

Fully parallel mapping of adaptive correlator algorithm. No time multiplexing.

$$N_{op} = 96$$

$$\text{Clock rate} = 25 \text{ MHz} \Rightarrow 2400 \text{ MOPS}$$

$$A_{op} = 5.4 \text{ mm}^2/96 = .15 \text{ mm}^2$$

$$\text{Power} = 12 \text{ mW}$$

# The Basic Problem is Time Multiplexing

---

- Processor architectures obtain performance by increasing the clock rate, because the parallelism is low
- Results in ever increasing memory on the chip, high control overhead and fast area consuming logic

*But doesn't time multiplexing give better area efficiency???*

# Area Efficiency

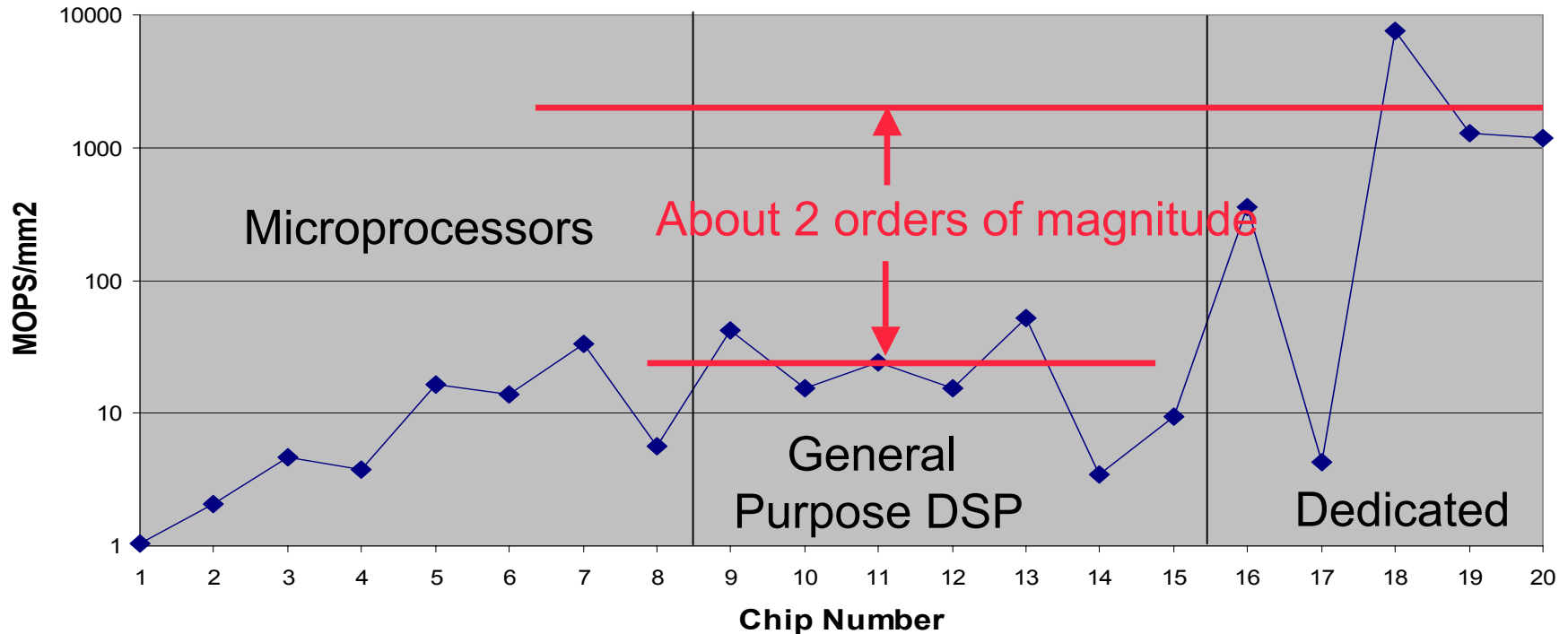
- SOC based devices are often very cost sensitive
- So we need a \$ cost metric => for SOC's it is equivalent to the efficiency of area utilization
- Area Efficiency Metric:

Computation per unit area = MOPS/mm<sup>2</sup>

*How much of a \$ cost (area) penalty will we have if we put down many parallel hardware units and have limited time multiplexing?*



# Surprisingly the area efficiency roughly tracks the energy efficiency...



The overhead of flexibility in processor architectures is so high that there is even an area penalty

# Hardware/software

---

Conclusion:

*There is no software/hardware tradeoff.*

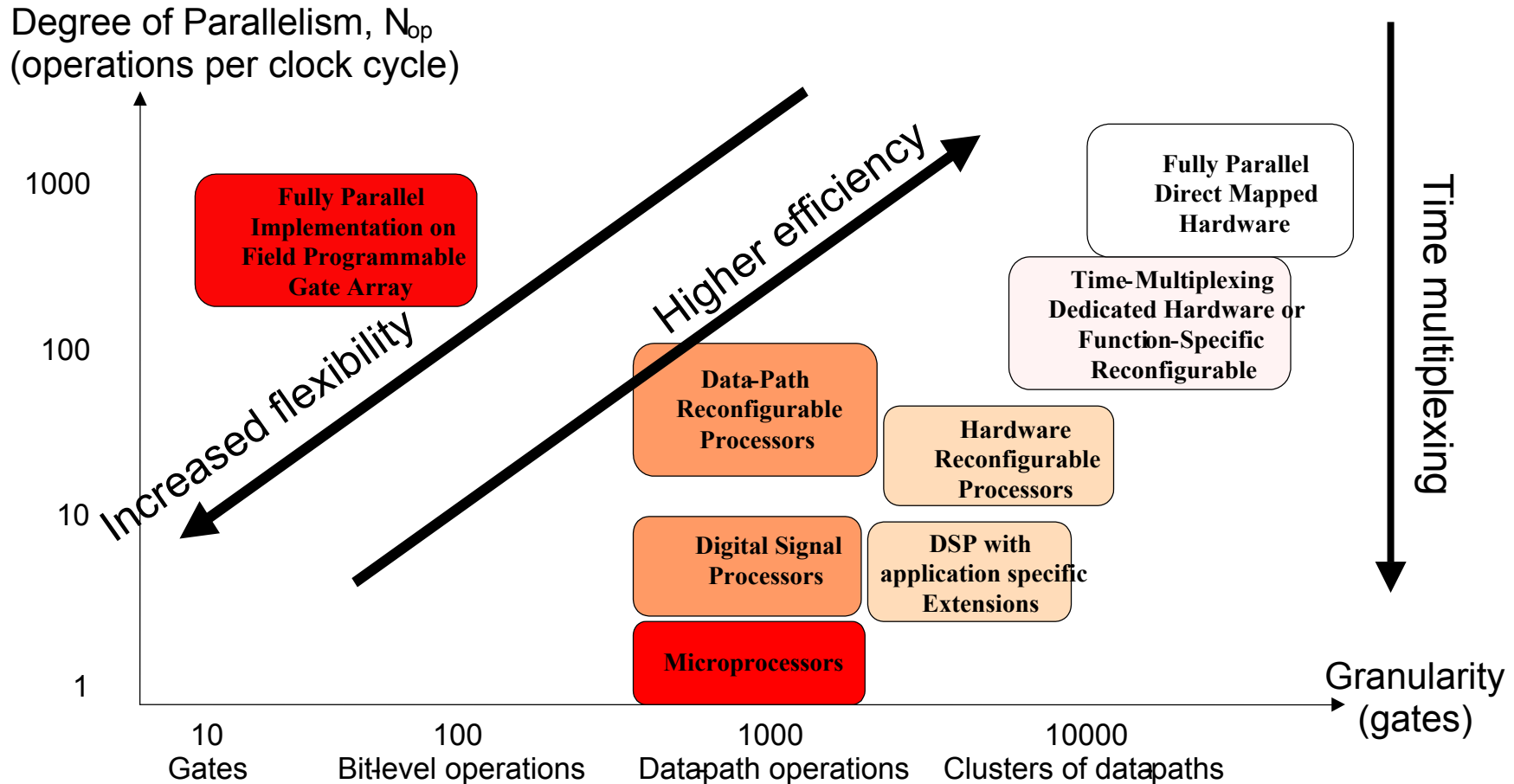
- The difference between hardware and software in performance, power and area is so large that there is no “tradeoff”.
- It is reasons other than power, energy, performance or cost that drives a software solution (e.g. business, legacy, ...).
- The **“Cost of Flexibility” is extremely high**, so the other reasons better be good!

# Are there better ways to provide flexibility?

---

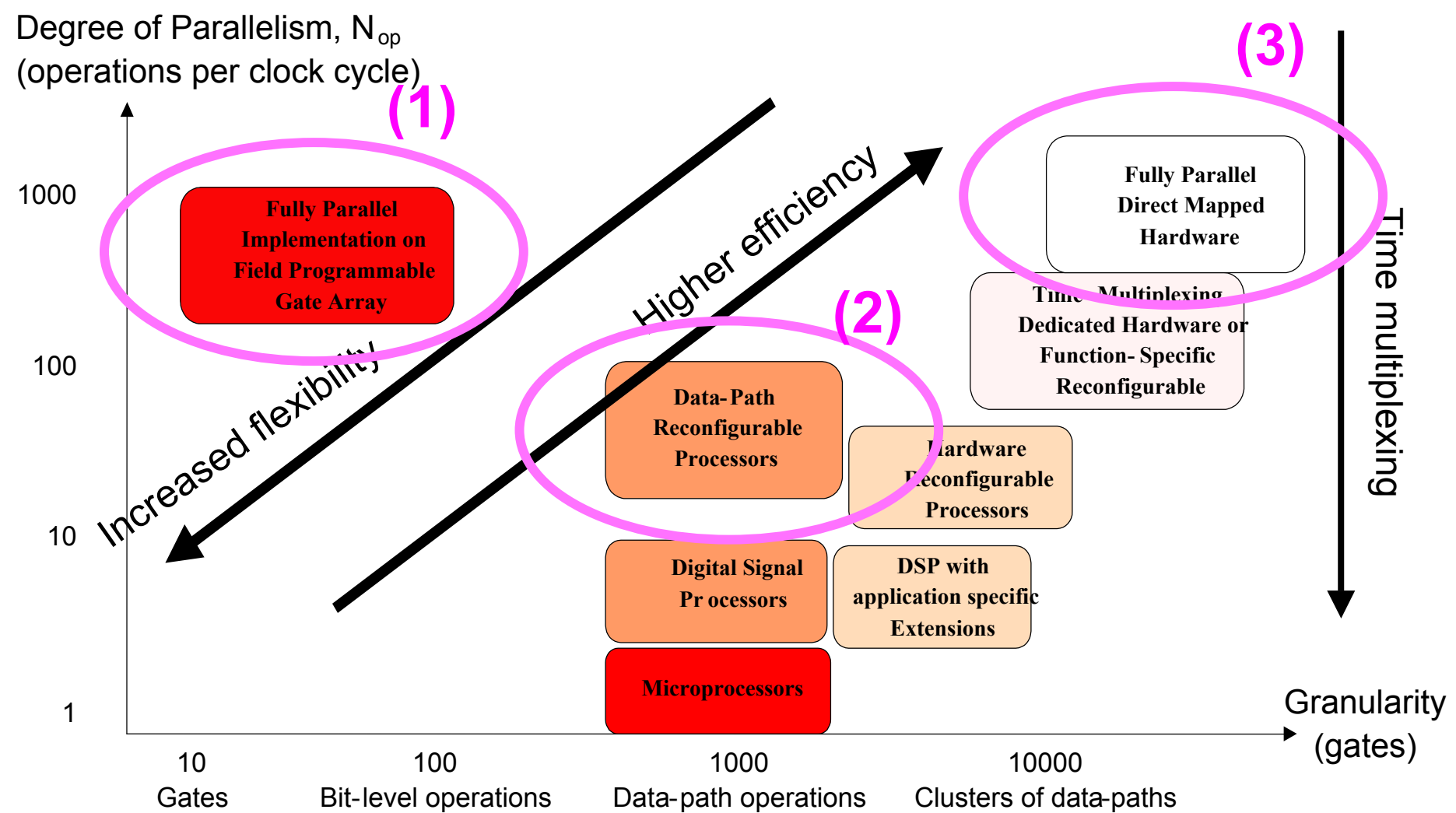
- Lets say the reasons for flexibility are good enough, then are there alternatives to processor based software programmability??
- Yes...
  - » The key is to provide flexibility along with the parallelism we get from the technology..
  - » Lots of choices...

# Granularity and Parallelism

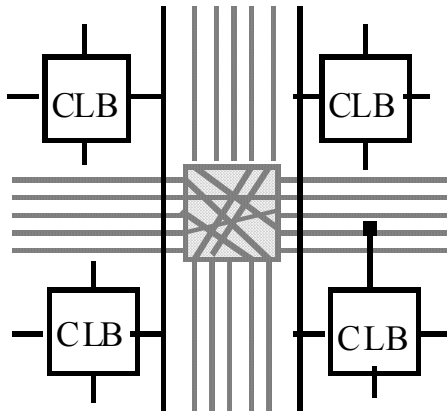


- Increased granularity and higher parallelism yields higher efficiency
- Smaller granularity and reduced parallelism yields more flexibility
- Time multiplexing is needed for performance with low parallelism

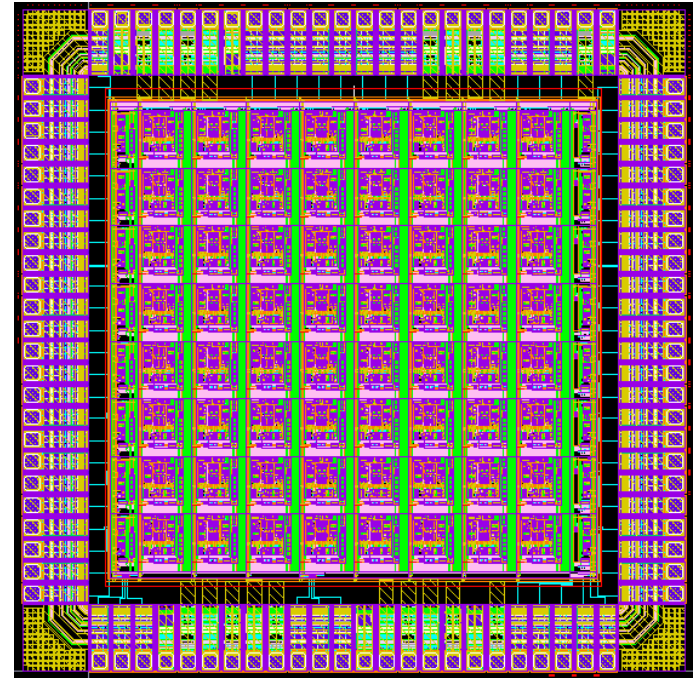
# We will look at three cases...



# Case (1): Reconfigurable Logic: FPGA

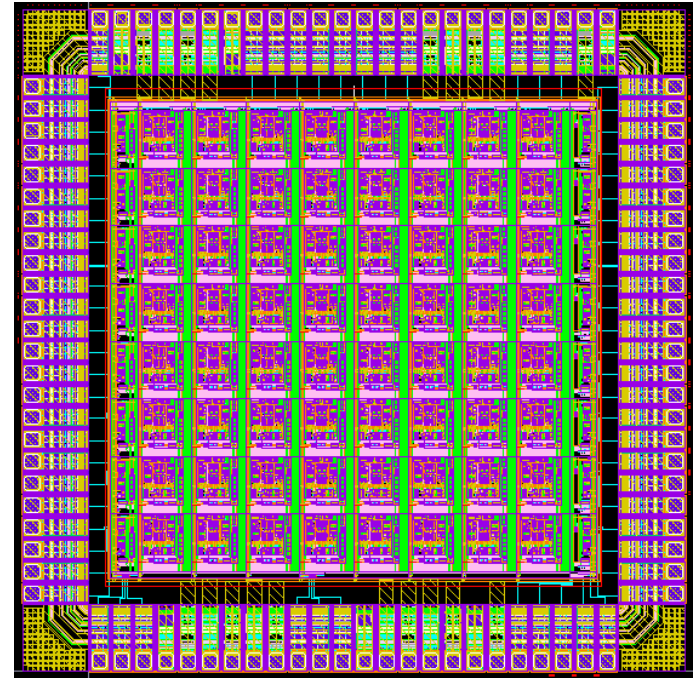
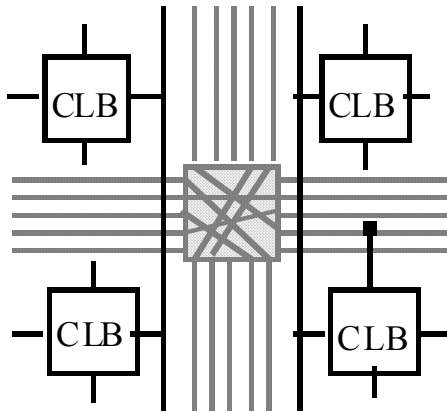


- Very low granularity (CLB's) – *improves flexibility*
- High parallelism – *improves efficiency*



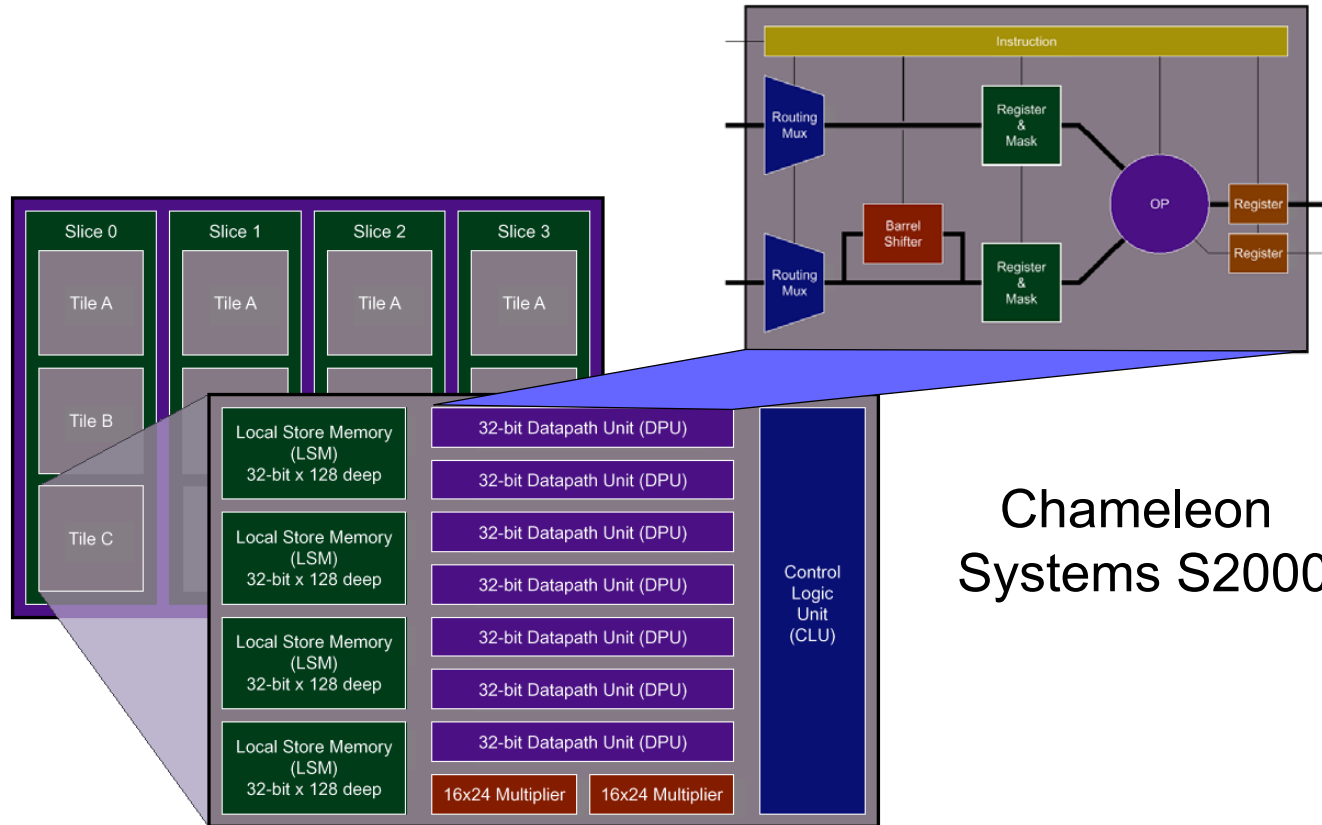
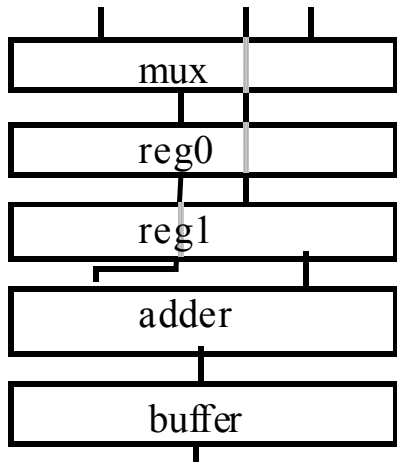
But....

# Case (1): Reconfigurable Logic: FPGA



- Very low granularity (high amount of interconnect) – decreases efficiency

# Case (2): Reconfiguration at a higher level of granularity

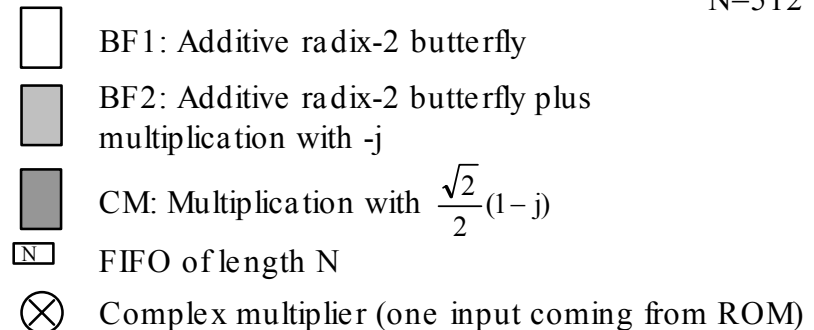
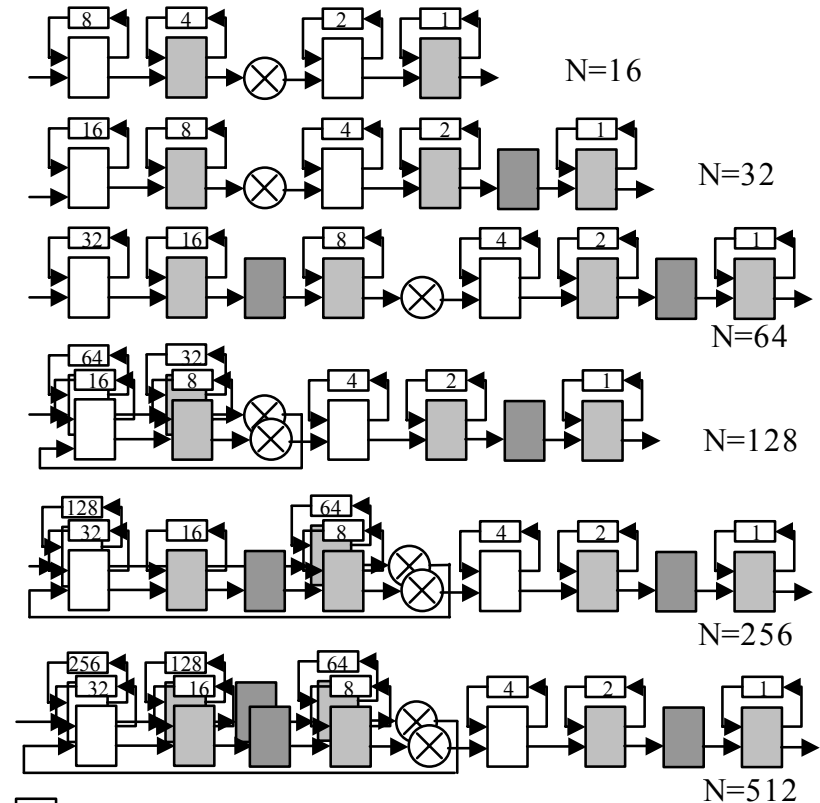


- Higher granularity – datapath units
- Higher efficiency, but lower flexibility

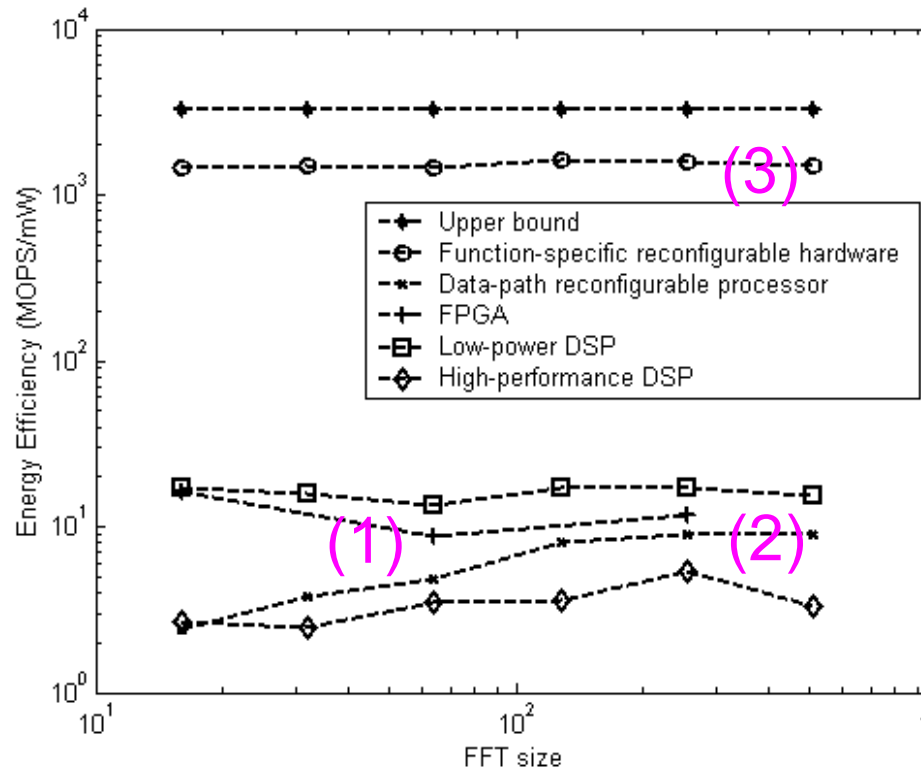


# Case (3): Even higher granularity - “Flexible” dedicated hardware

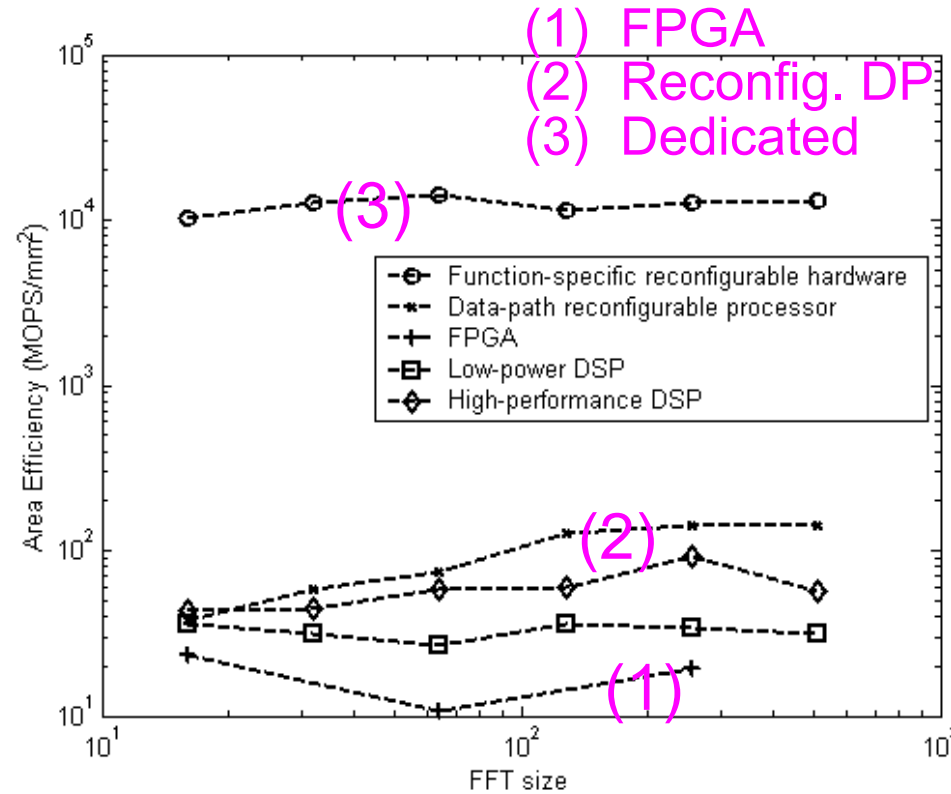
- Use a hardware architecture that has the flexibility to cover a range of parameter values
- Not much flexibility, but very high efficiency
- Example here is an FFT which can range from N=16 to 512
- Uses time multiplexing



# Efficiencies for a variety of architectures for a flexible FFT



MOPS/mW  
vs. FFT size



MOPS per mm<sup>2</sup>  
vs. FFT size

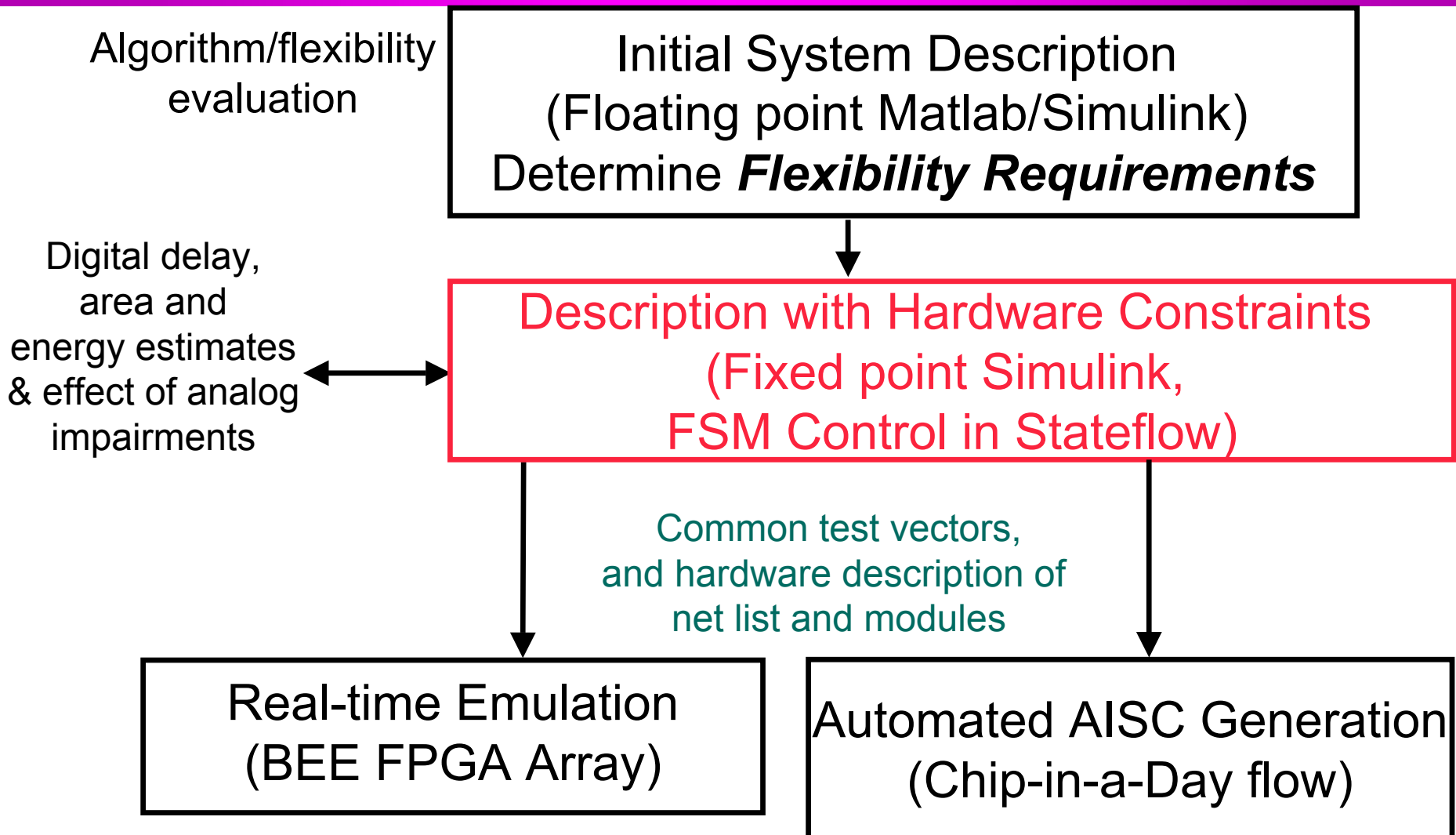
\* All results are scaled to 0.18 $\mu$ m

# The Issues

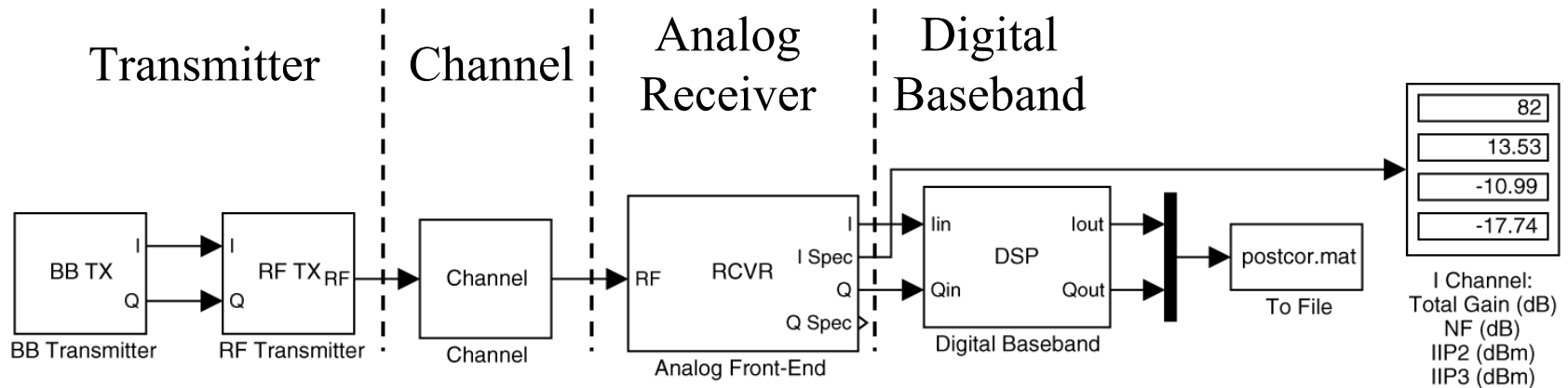
---

- How much flexibility is needed and how best to include it...
- A single system description including interaction between the analog and digital domains
- “Realtime” SOC prototyping
- Automated ASIC design flow

# An SOC Design Flow with Prototyping

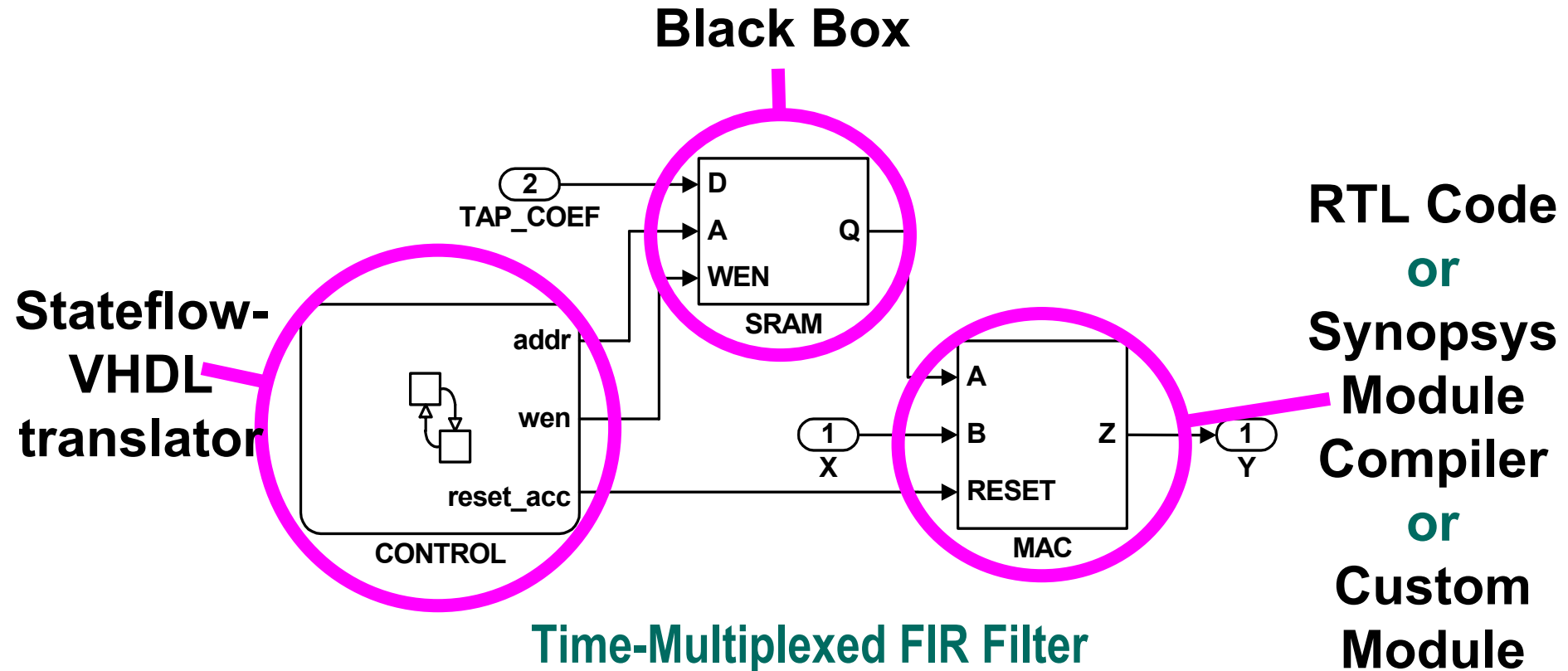


# Simulation Framework using Simulink/Stateflow (from Mathworks, Inc.)



- **Techniques used to decrease simulation time:**
  - Baseband-equivalent modeling of RF blocks
  - Compile design using MATLAB Real-Time Workshop

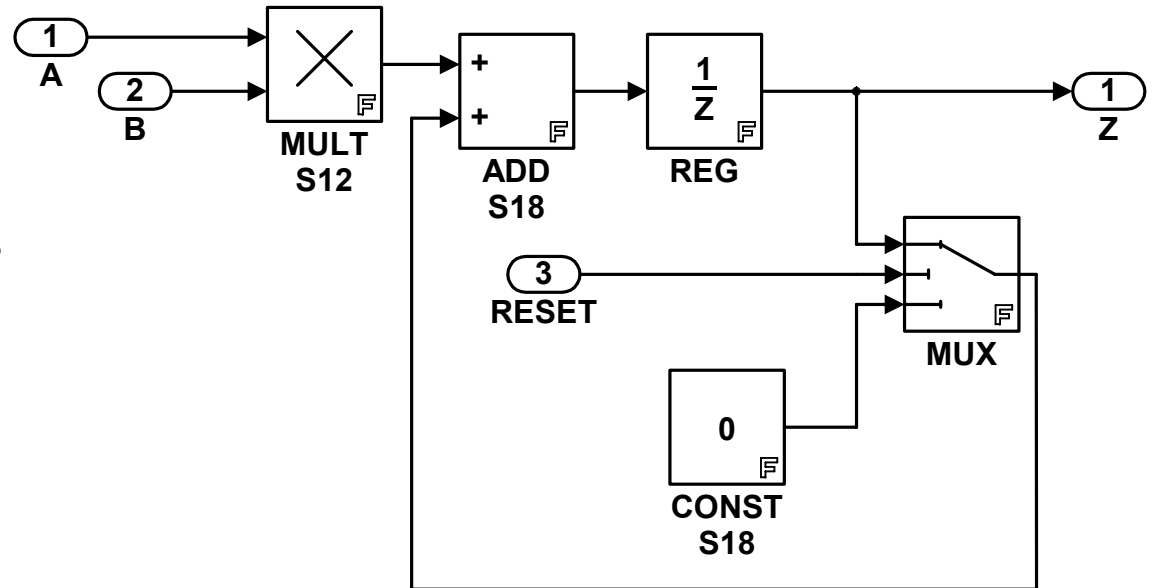
# Blocks map to implementation libraries



- Implementation choices embedded in description
- Libraries of blocks are pre-verified and re-used

# Timed Dataflow Graph Specification

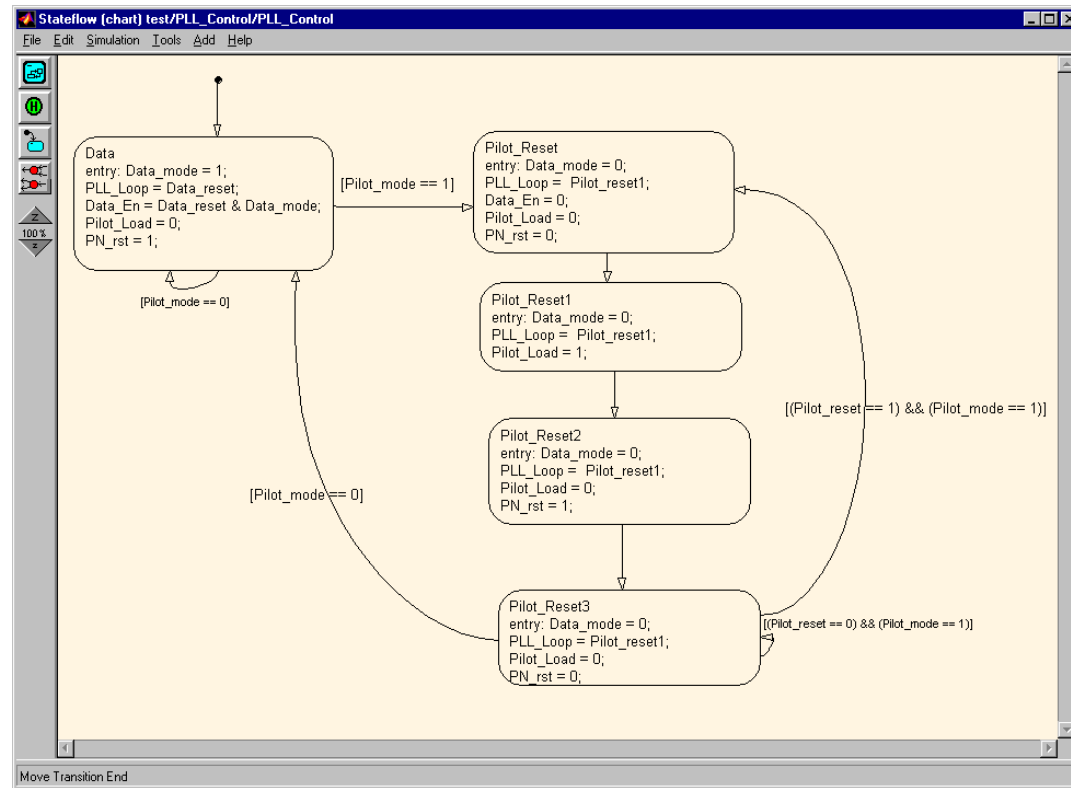
- Simulink (from Mathworks)
- Discrete-Time (cycle accurate)
- Fixed-Point Types (bit true)
- No need for RTL simulation
- Embedded implementation choices



Multiply / Accumulate

# Control

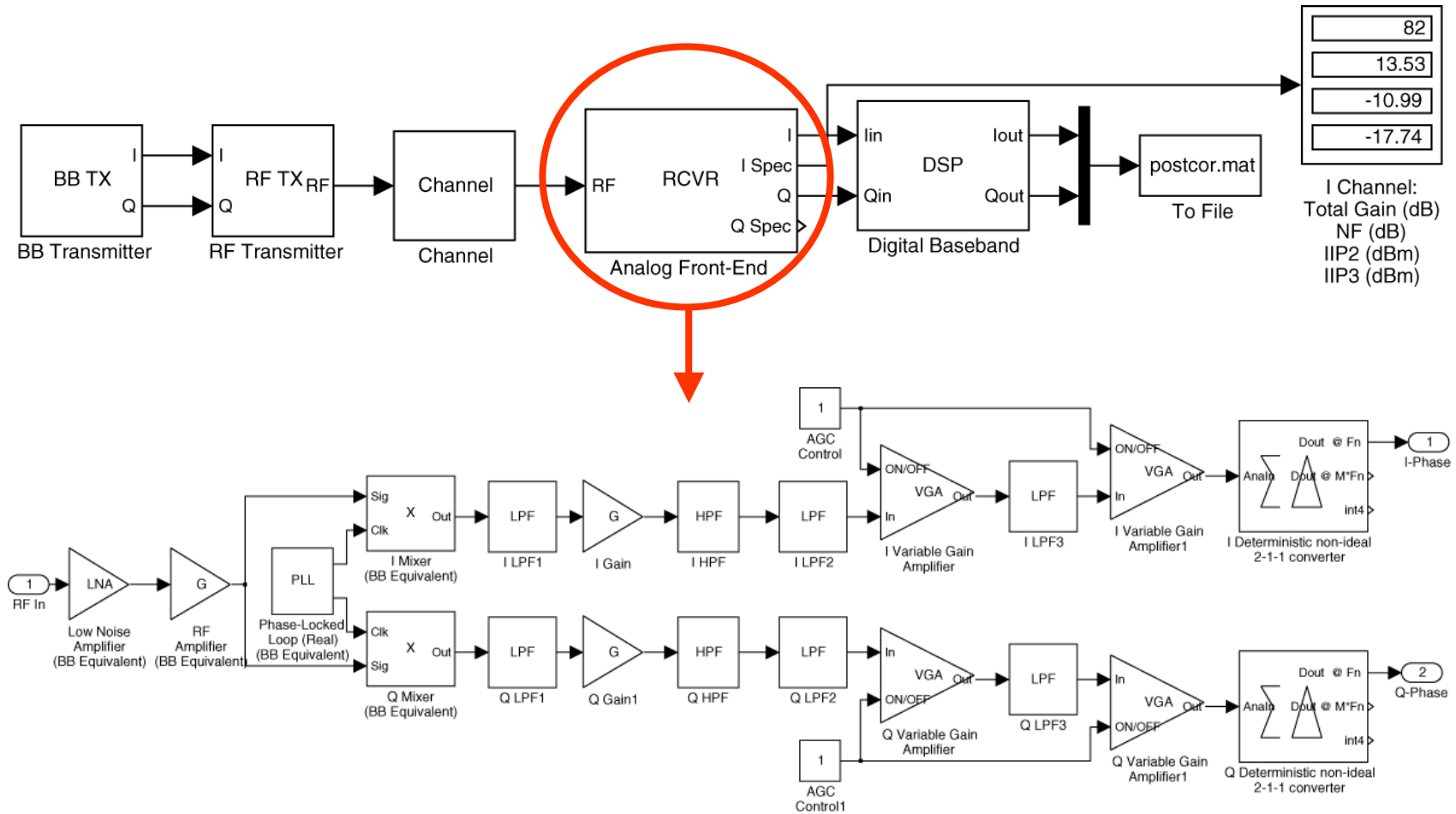
- Stateflow
  - » Extended Finite State Machine
  - » Subset of Syntax
  - » Converted to VHDL
  - » Synthesized
- VHDL
  - » Synthesized directly



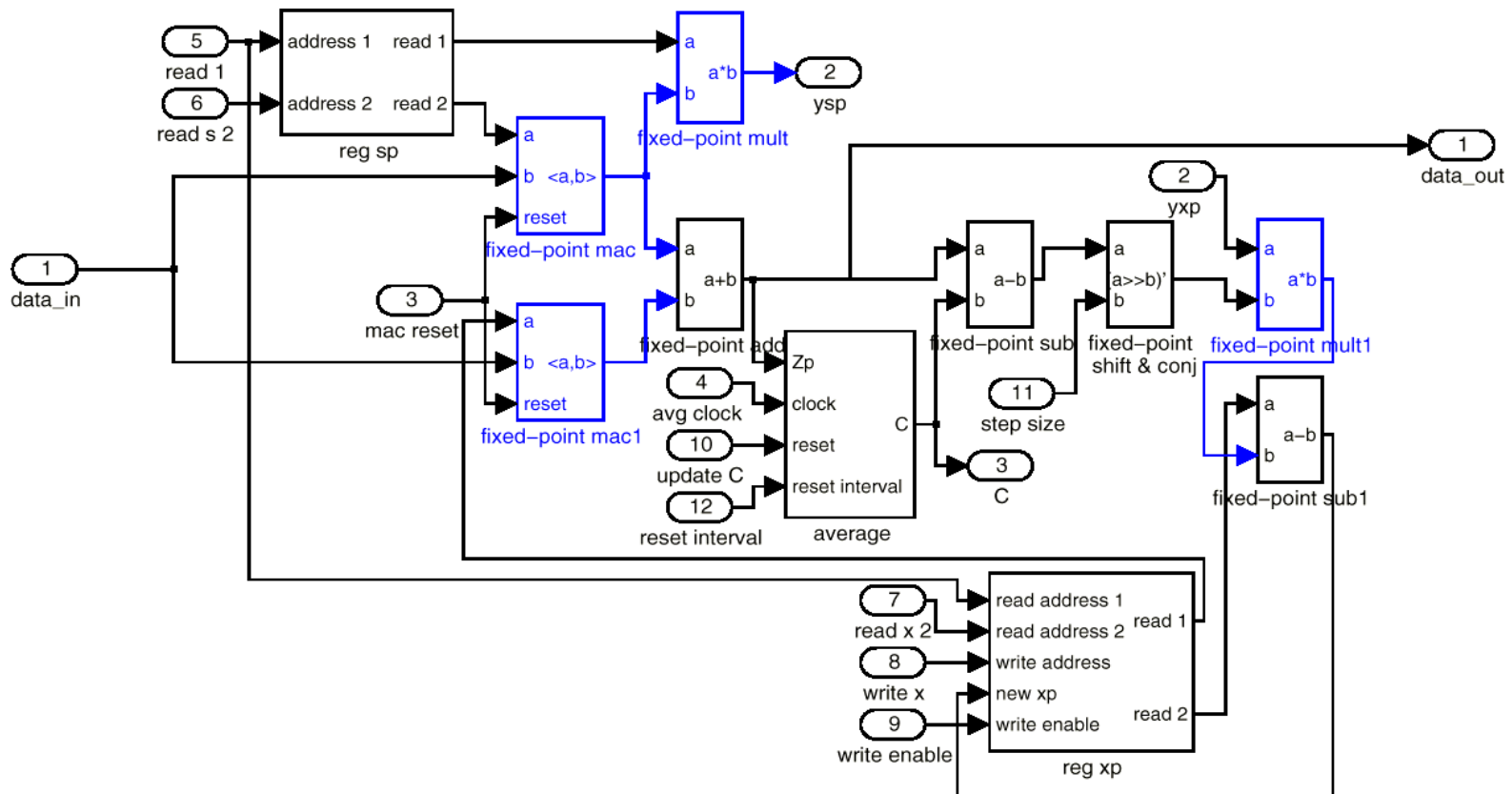
VHDL & Stateflow Macros map to a netlist of Standard Cells using standard synthesis



# Simulink Model of Direct-Conversion Receiver



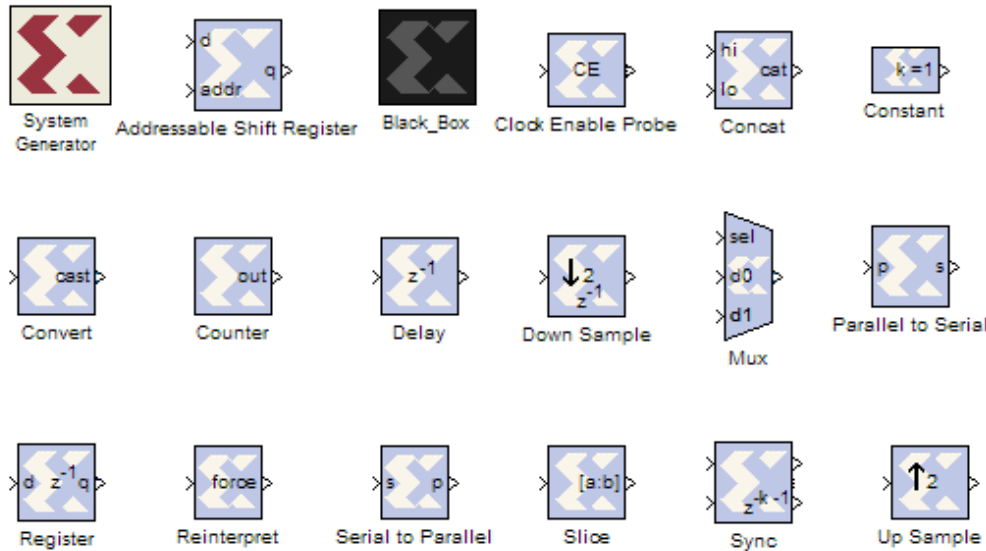
# Bit true, cycle accurate digital baseband algorithms...



# Basic Blocks based on Xilinx System Generator libraries

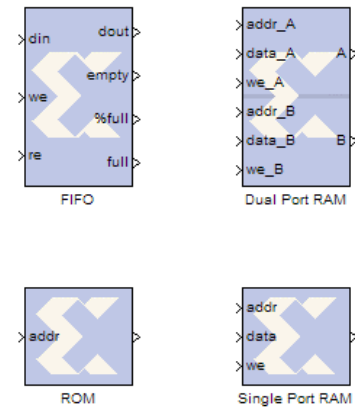
Xilinx Blockset v2.2  
(c) 2002 Xilinx, Inc.

## Basic Library



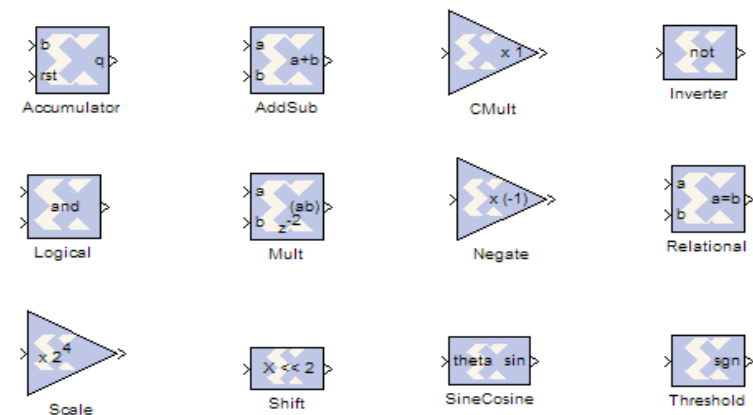
Xilinx Blockset v2.2  
(c) 2002 Xilinx, Inc.

## Memory Library



Xilinx Blockset v2.2  
(c) 2002 Xilinx, Inc.

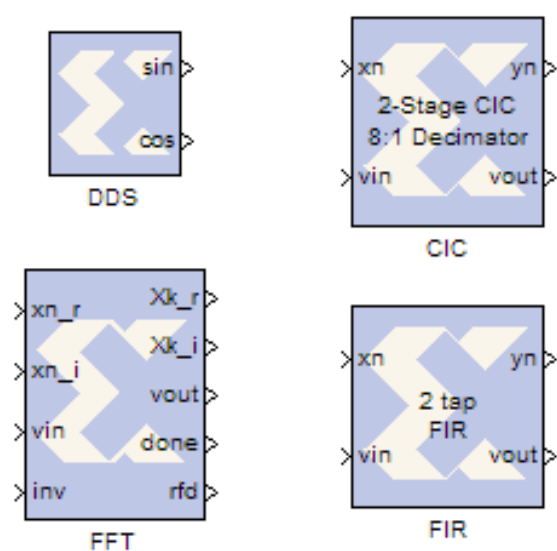
## Math Library



# Higher level DSP Blocks

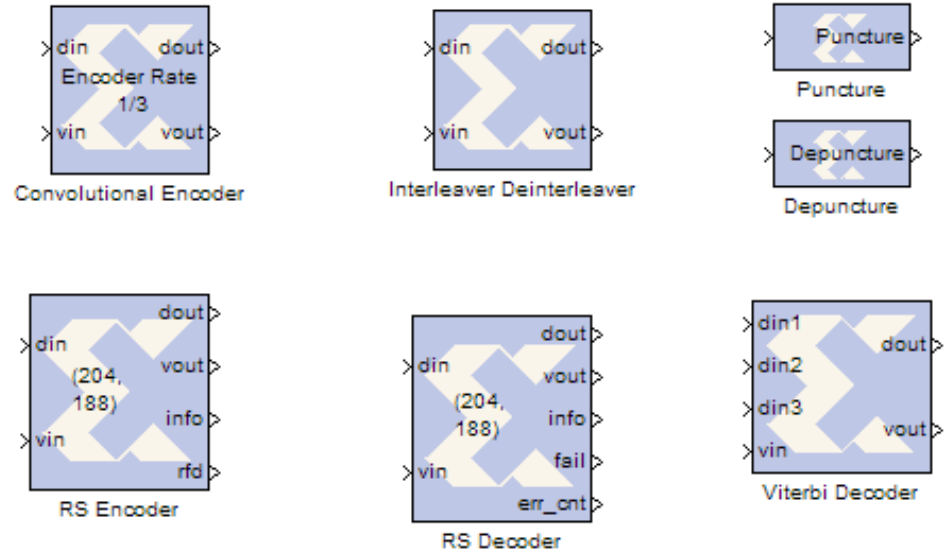
Xilinx Blockset v2.2  
(c) 2002 Xilinx, Inc.

DSP Library

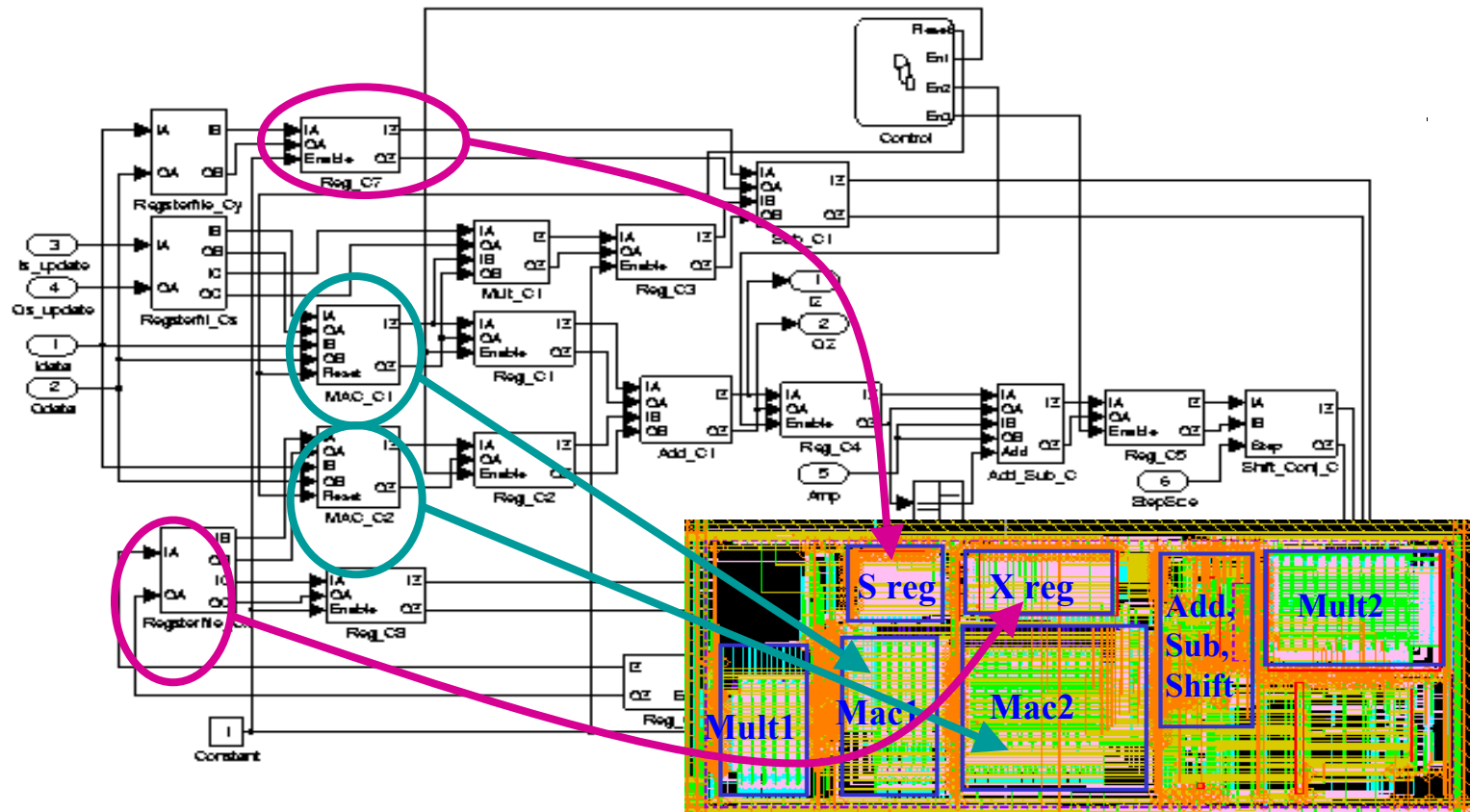


Xilinx Blockset v2.2  
(c) 2002 Xilinx, Inc.

Communication Library



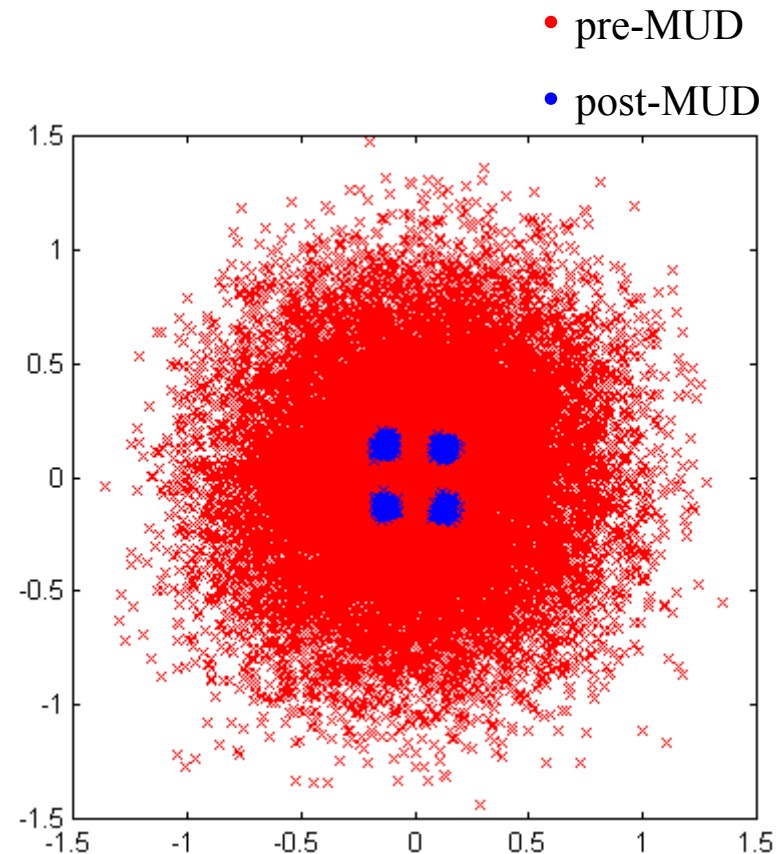
Directly map diagram into hardware since there is a one for one relationship for each of the blocks



- Results: A fully parallel architecture that can be implemented rapidly

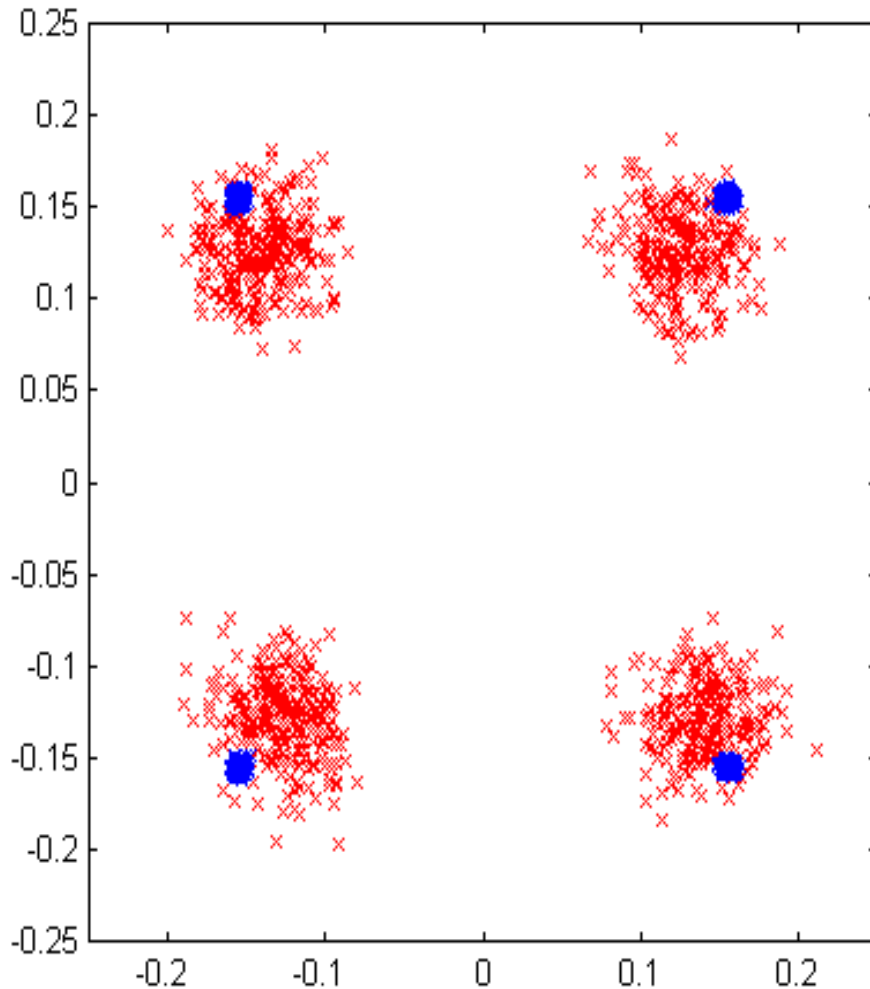
# Then do a simulation: Zero-IF Receiver

- 10 users (equal power)
- 13.5dB receiver NF
- PLL: -80dBc/Hz @ 100kHz
- 2.5° I/Q phase mismatch
- 82dB gain
- 4% gain mismatch
- IIP2 = -11dBm
- IIP3 = -18dBm
- 500kHz DC notch filter
- 20MHz Butterworth LPF
- 10-bit, 200MHz  $\Sigma$ - $\Delta$  ADC



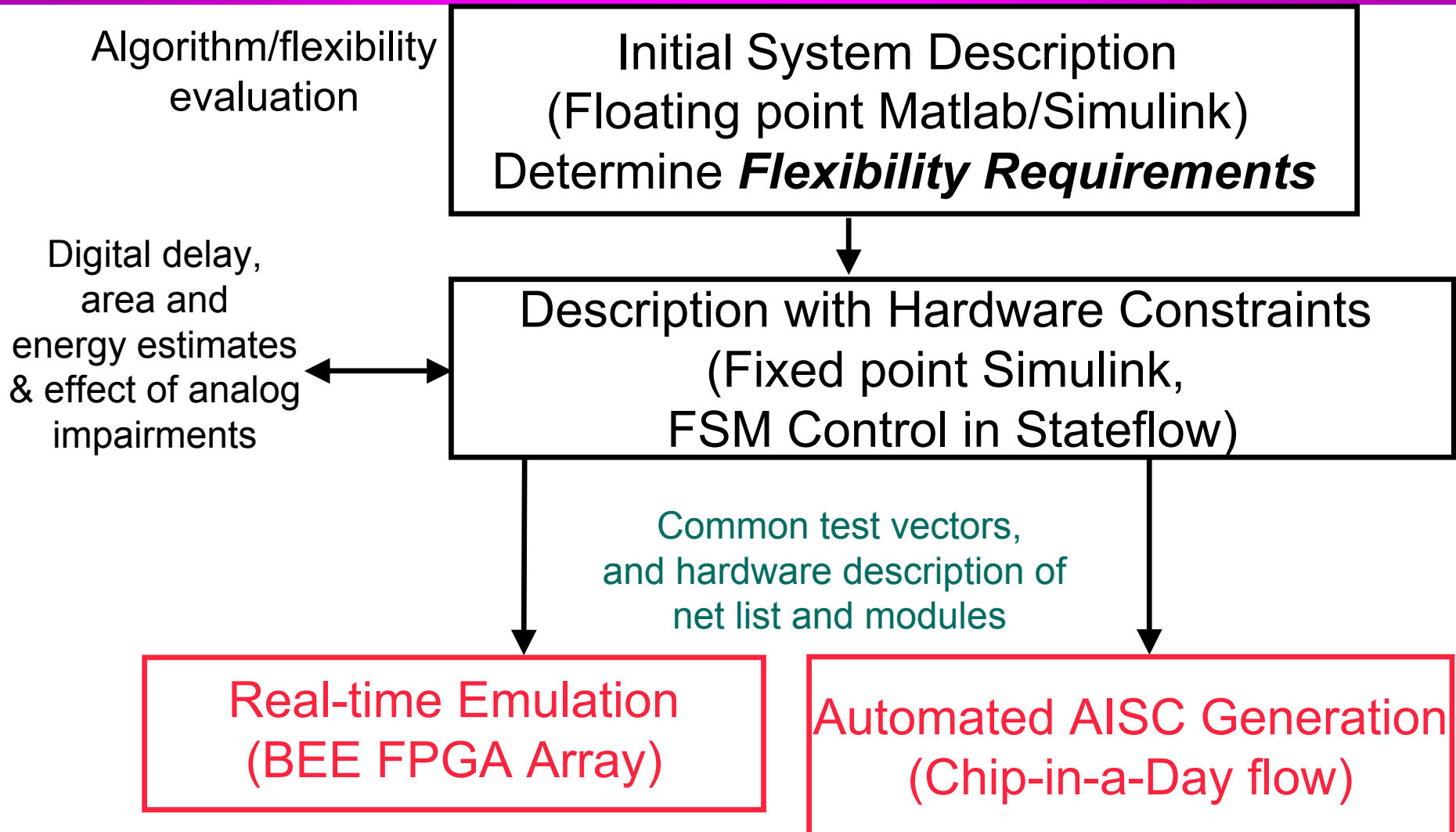
**Output SNR  $\approx$  15dB**

# With Analog Impairments



- 10 users (equal power)
- 20MHz Butterworth LPF
- 500kHz DC notch filter
- 13.5dB receiver NF
- 82dB gain
- 4% gain mismatch
- 2.5° I/Q phase mismatch
- IIP2 = -11dBm
- IIP3 = -18dBm
- PLL: -80dBc/Hz @ 100kHz
- 10-bit, 200MHz S-D ADC

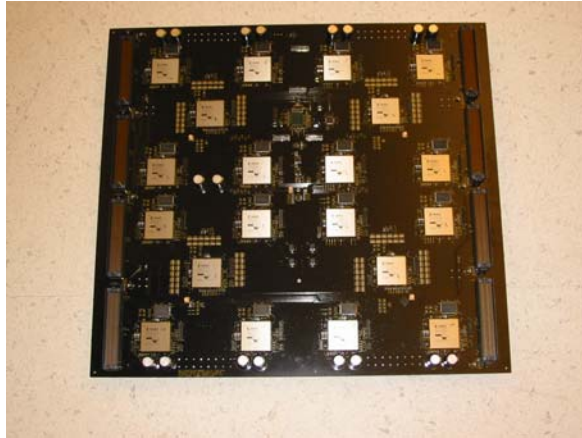
# Now to implement that description



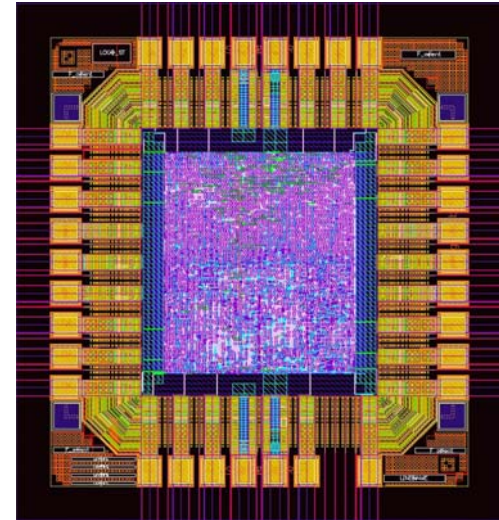


# Single description – Two targets

Simulink/Stateflow  
Description



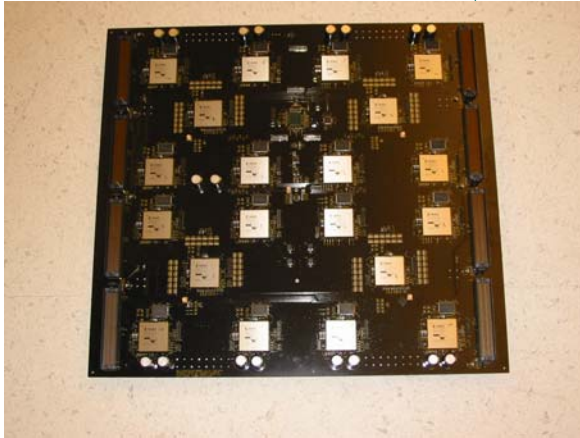
BEE  
FPGA Array



ASIC Implementation  
“Chip in a day”

# BEE Target for Real-time emulation

Simulink/Stateflow  
Description



BEE  
FPGA Array

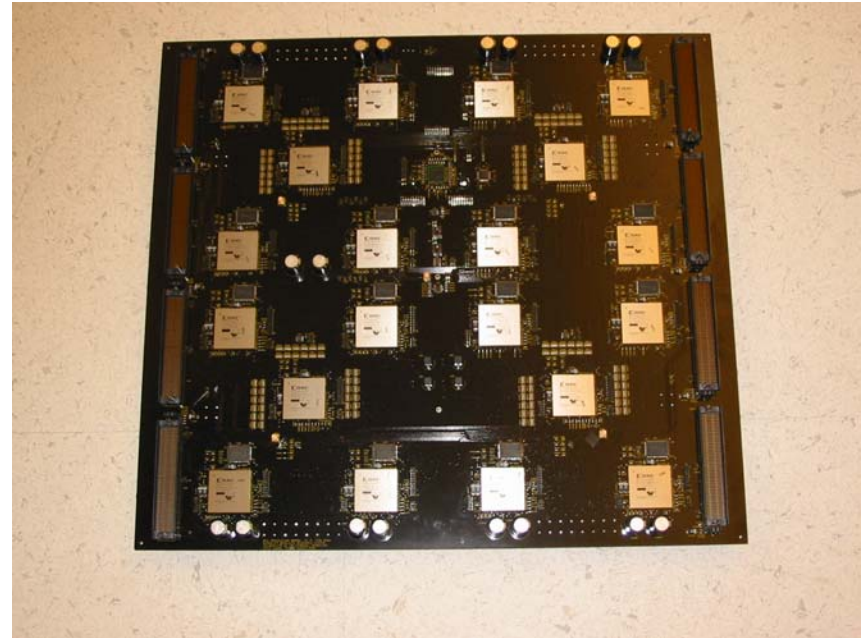
# BEE Design flow Goals

---

- Fully automatic generation of FPGA and ASIC implementations from Simulink system level design
- Cycle accurate bit-true functional level equivalency between ASIC & BEE implementation
- Real-time emulation controlled from workstation

# Processing Board PCB

- Board-level Main Clock Rate: 160MHz+
- On Board connection speed:
  - » FPGA to FPGA: 100MHz
  - » XBAR to XBAR: 70MHz
- Off board connection speed: (3 ft SCSI cable loop back through riser card)
  - » LVTTTL: 40MHz
  - » LVDS: 160MHz ~ 220MHz



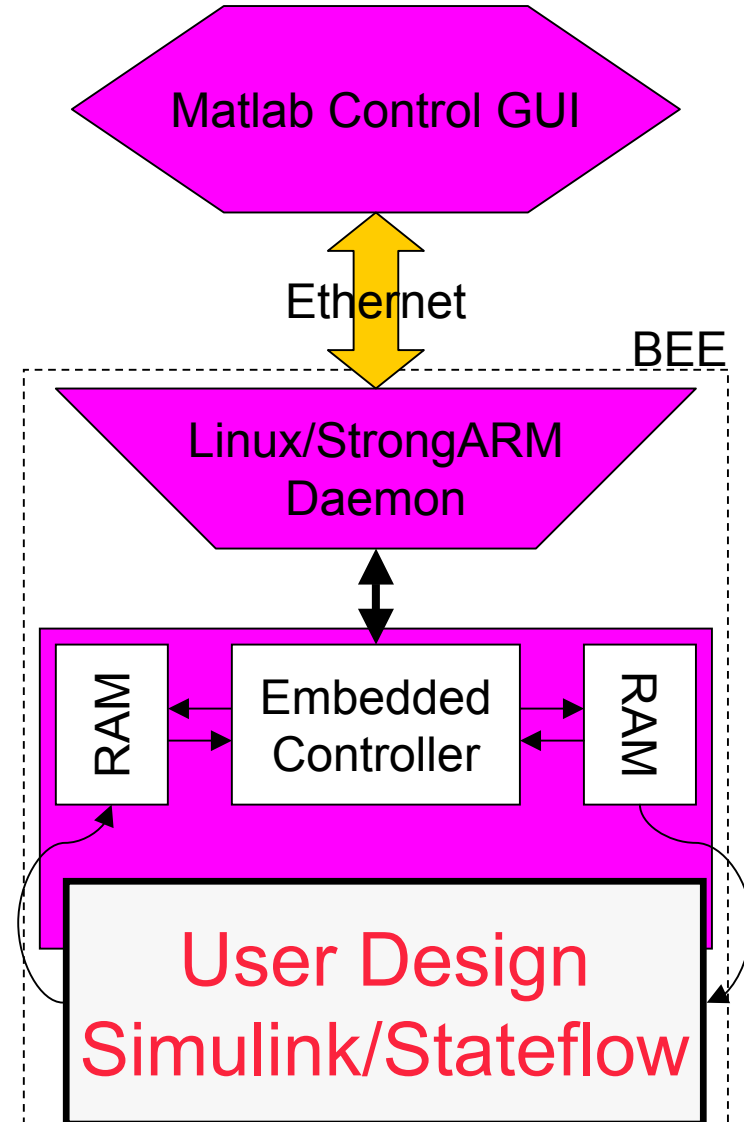
- Board Dimension: 53 X 58 cm
- Layout Area: 427 sq. in.
- No. of Layers: 26

# The BEE with RF transceiver I/O

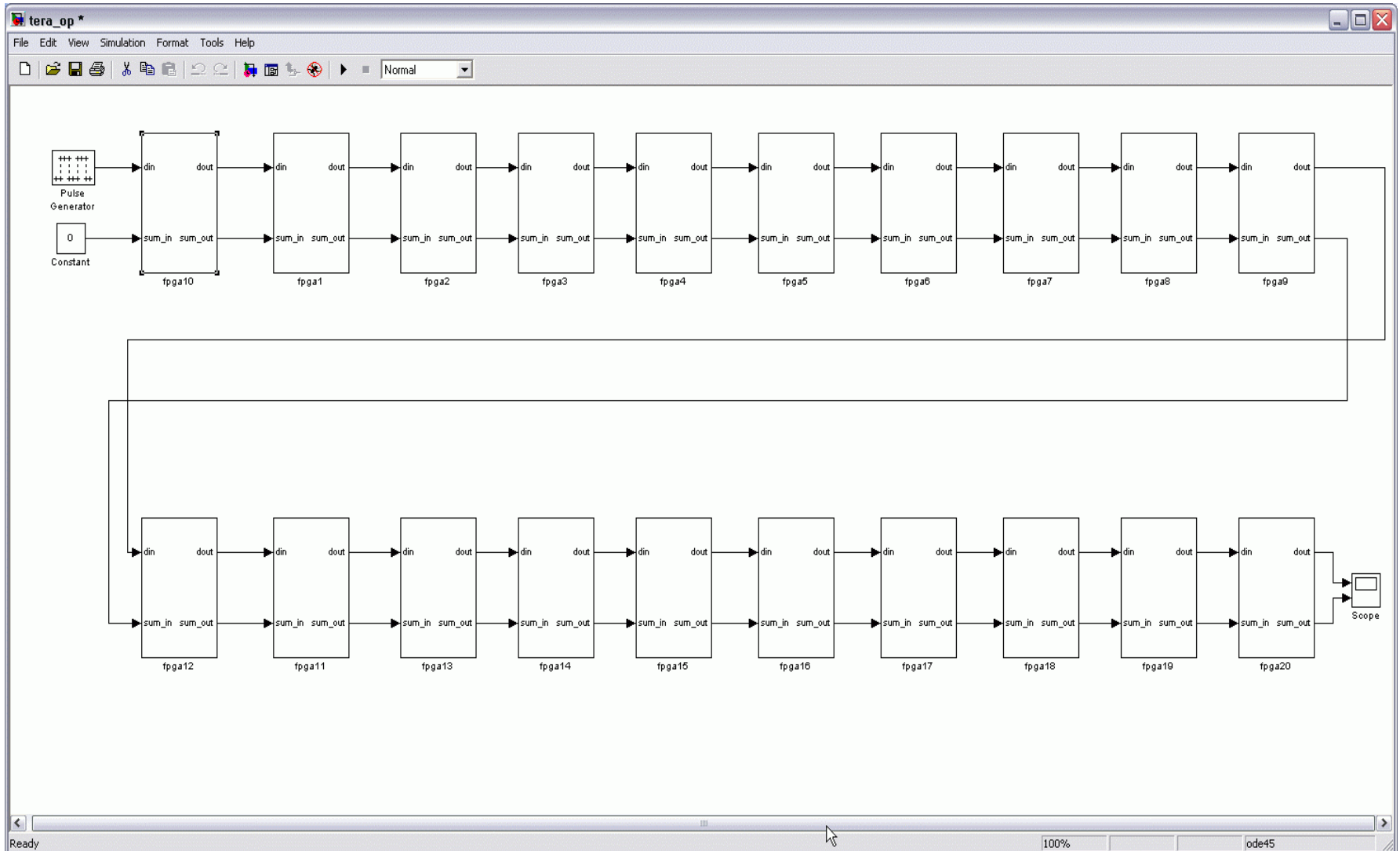


# Run-time Data I/O Interface

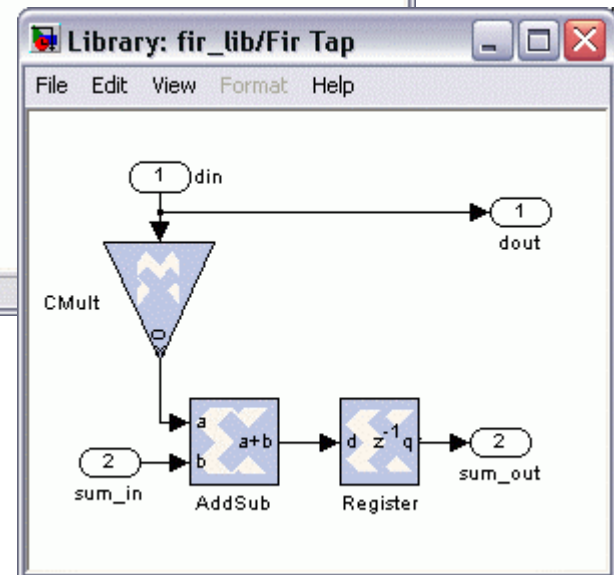
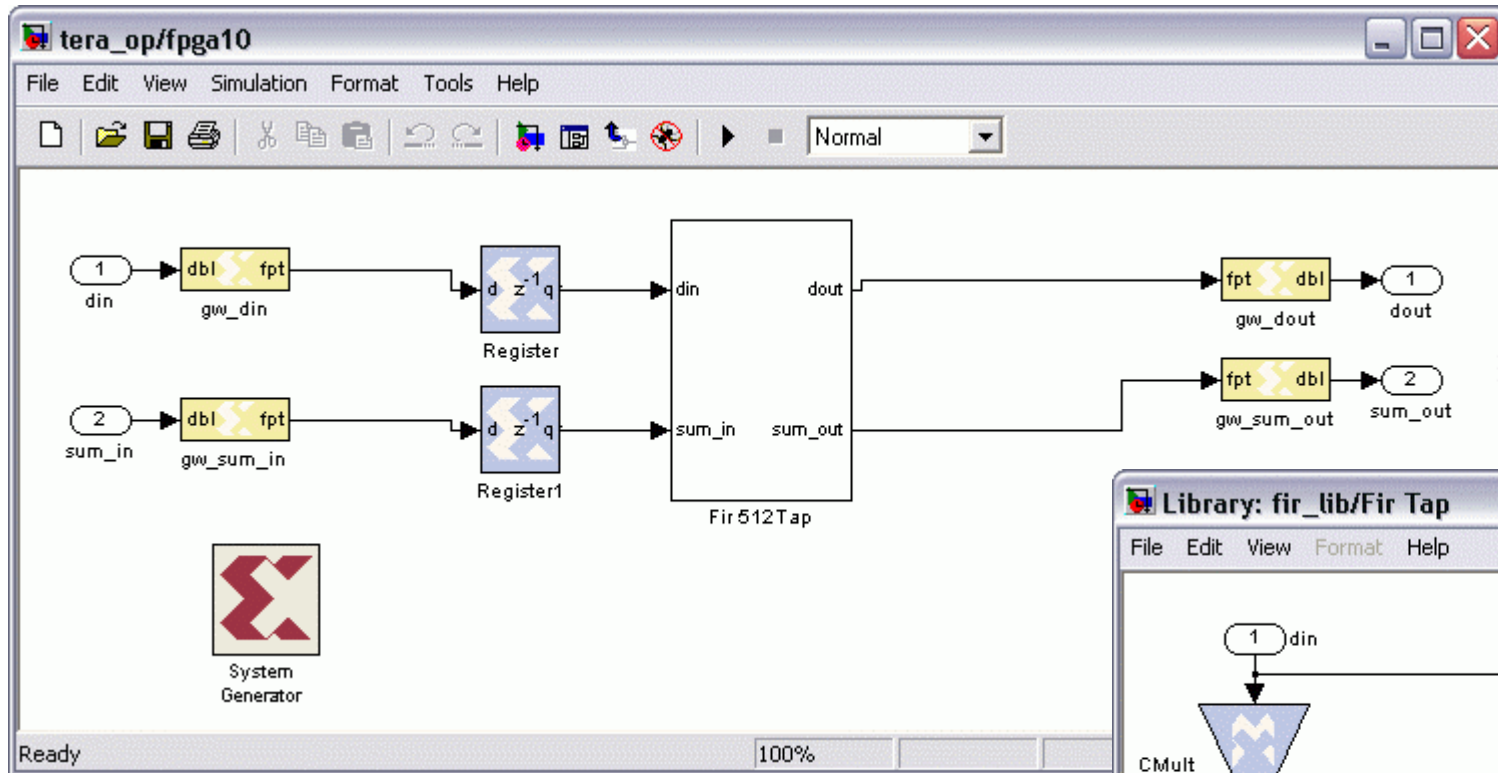
- Infrastructure for transferring data to and from the BEE
  - » The entire hardware interface is in one fully parameterized block
  - » Simply drop block into the Simulink diagram
  - » Accepts standard Simulink data structures for reuse of existing test vectors



# Benchmark: 10240 Tap Fir Design



# 10240 Tap Fir Design (cont.)





# BEE Performance

- Reference Design:
  - » 10240 tap FIR filter
  - » 512 taps per FPGA
- Slice utilization: 99% of 19200 slices
- Max Clock Rate: 30 Hz
- MOPS: 580,000 MOPS total (16bit add & 12bit cmult)
- Power: 2.5W per FPGA, 50W total

Comparison with an ASIC version using .13 micron chip metrics of 5000 MOPS/mm<sup>2</sup>, 1000 MOPS/mW =>

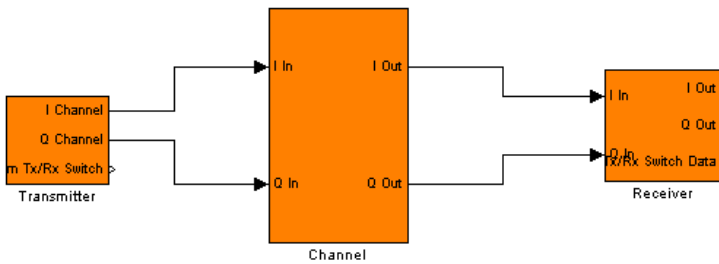
*The BEE is equivalent to a single chip of 50 mm<sup>2</sup> with power = 500 mW.*

50 Watts/500 mW => 100 times more power

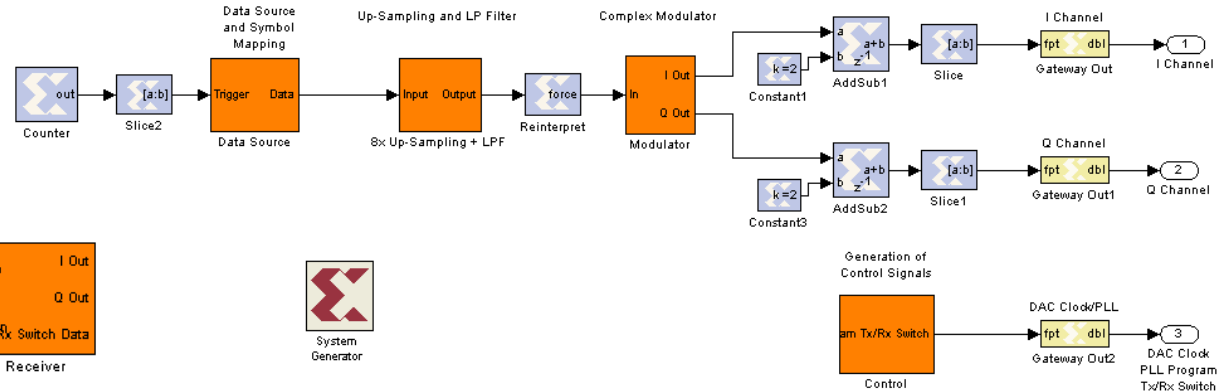
(20 \* 2 cm<sup>2</sup>)/.5 => 100 times more area

# Implementation of a Narrow-Band Radio System (Hans Bluethgen)

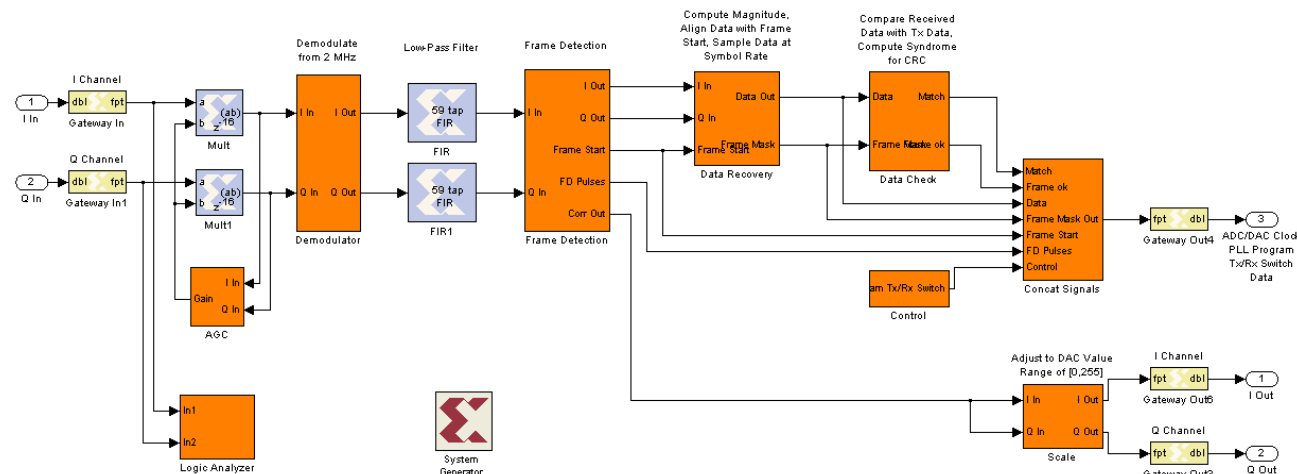
Complete System



Transmitter



Receiver



Data Rate 1 Mbit/s, 500 Kbit/s

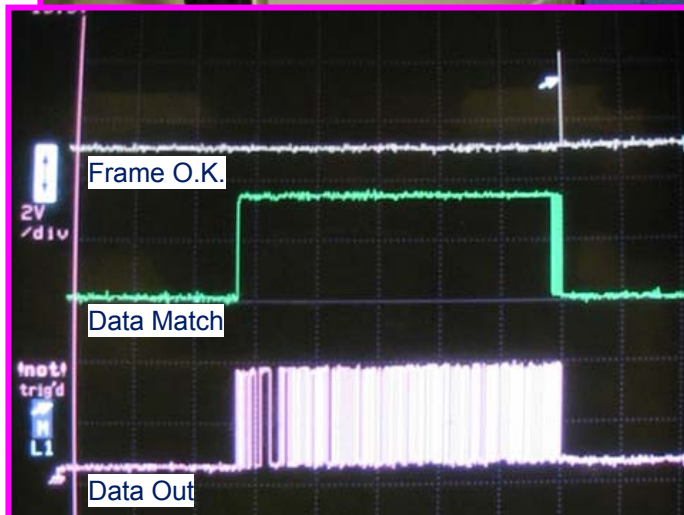
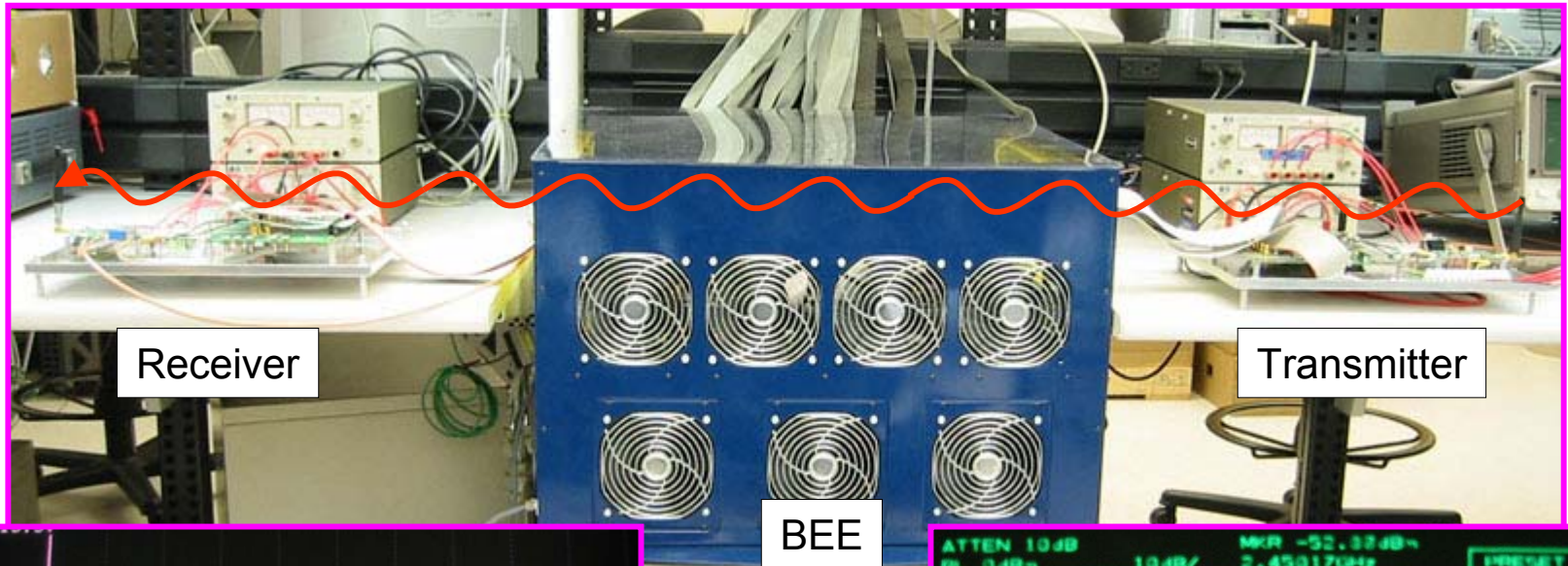
Carrier Freq. 2.45 GHz

Bandwidth 1 MHz

Modulation DPSK

Frame Synch. PN Sequence

# BEE Implementation of a Narrow-Band Radio

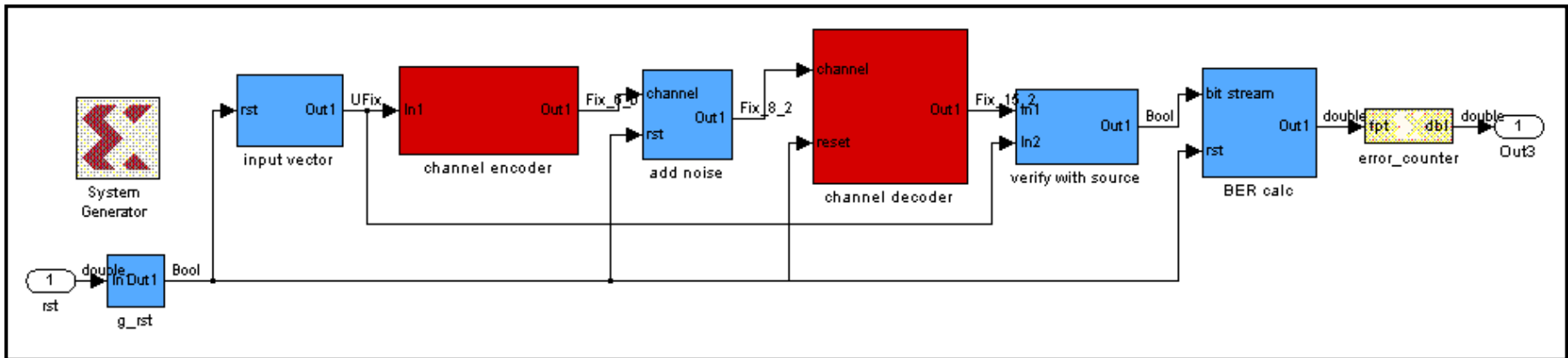


Receiver Output on SCSI Connector



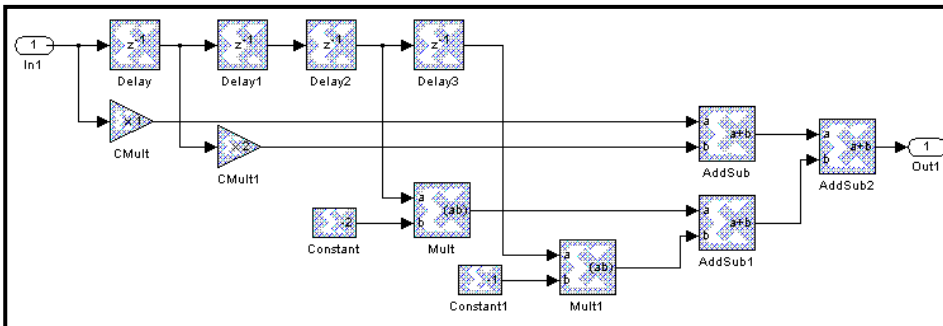
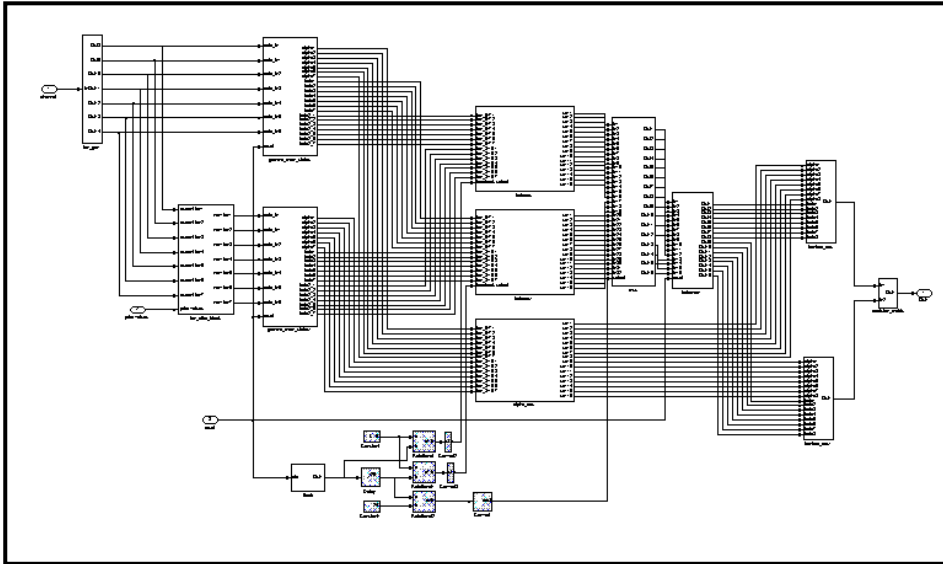
Transmitter Output Spectrum

# 3G Turbo Decoder (Bora Nikolic)



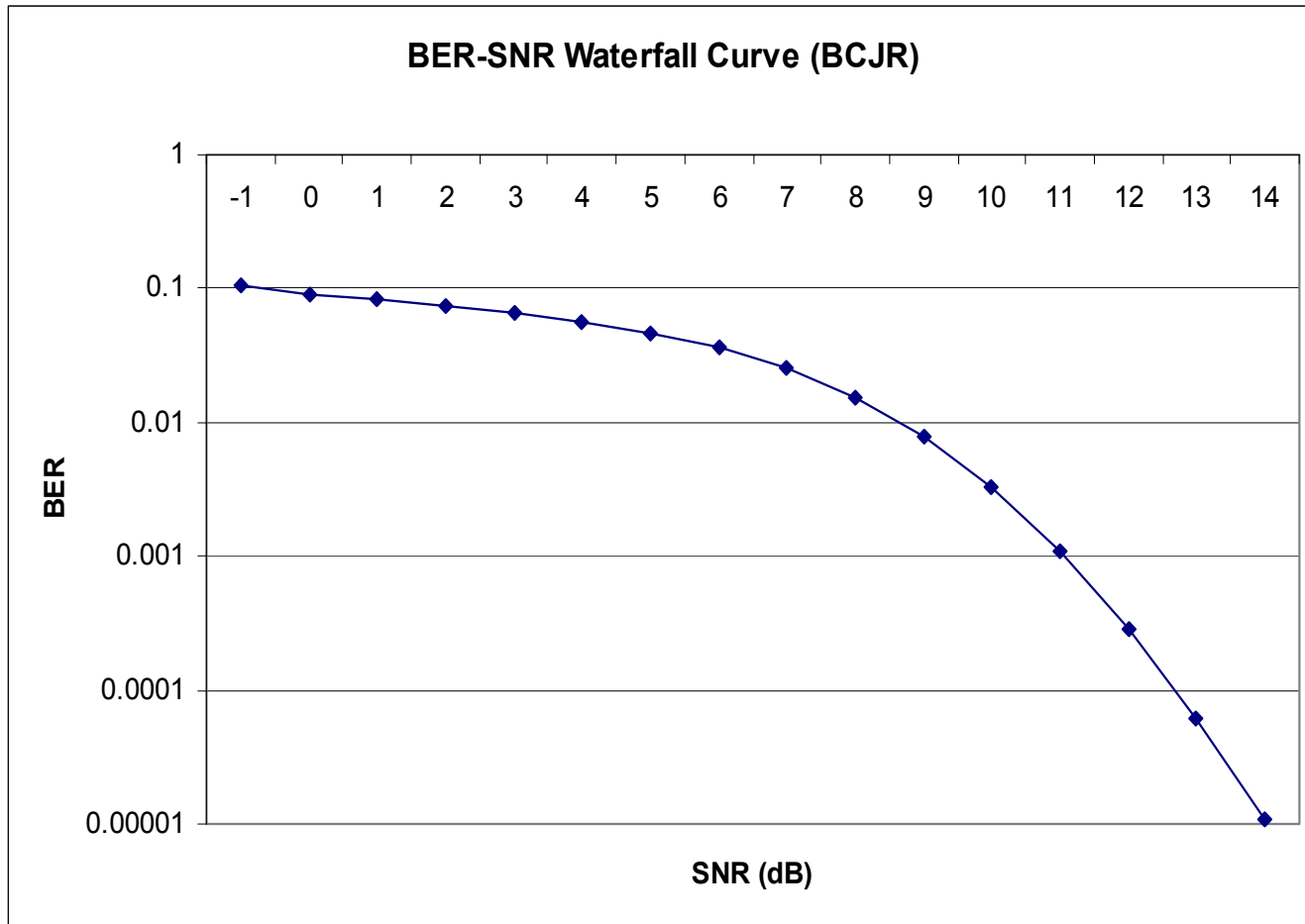
- Complete description of ECC with variable noise levels to evaluate performance
- 10 MHz system clock
- SNR 14db  $\rightarrow$  -1db
- $10^9$  Samples in two minutes
- Parameterized to support variable binary point precision, SNR, number of samples for architectural evaluation

# BCJR Simulink simulation



- E2PR4 Channel Encoder - Decoder
- Fully enclosed design
  - » Uniform RNG input vector
  - » Channel encoder
  - » AWGN filter
  - » Channel decoder
  - » BER collection mechanism
- Part of: Full 3G Turbo Decoder

# BCJR Waterfall Curve

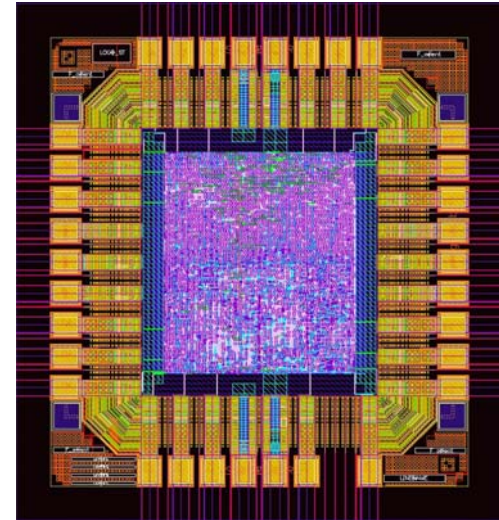


10MHz,  $10^9$  Samples, 1 bit binary point precision

Total simulation: approx. 10 minutes

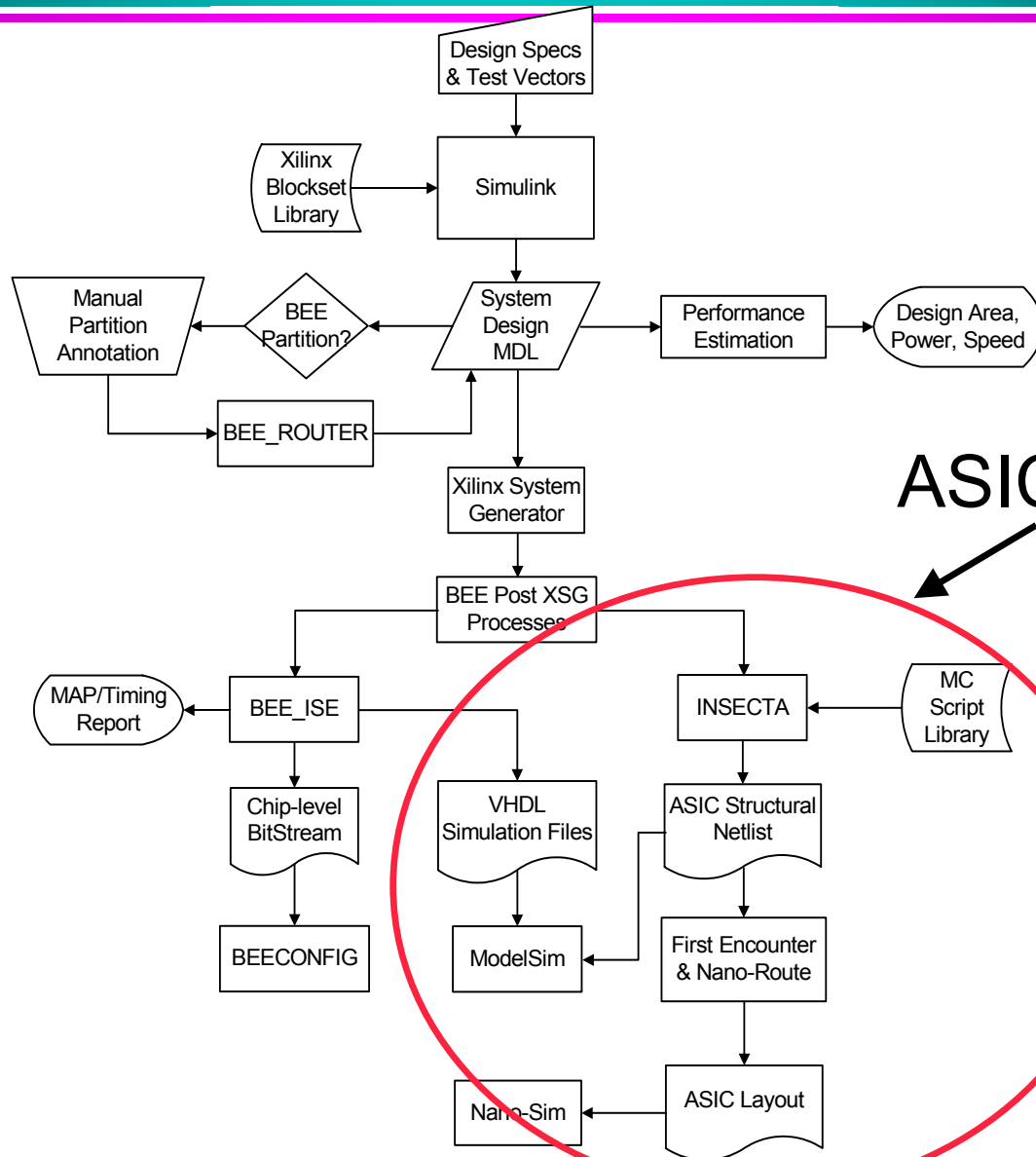
# ASIC Target

Simulink/Stateflow  
Description



ASIC Implementation  
“Chip in a day”

# Complete Design Flow





# Chip-in-a-Day ASIC flow

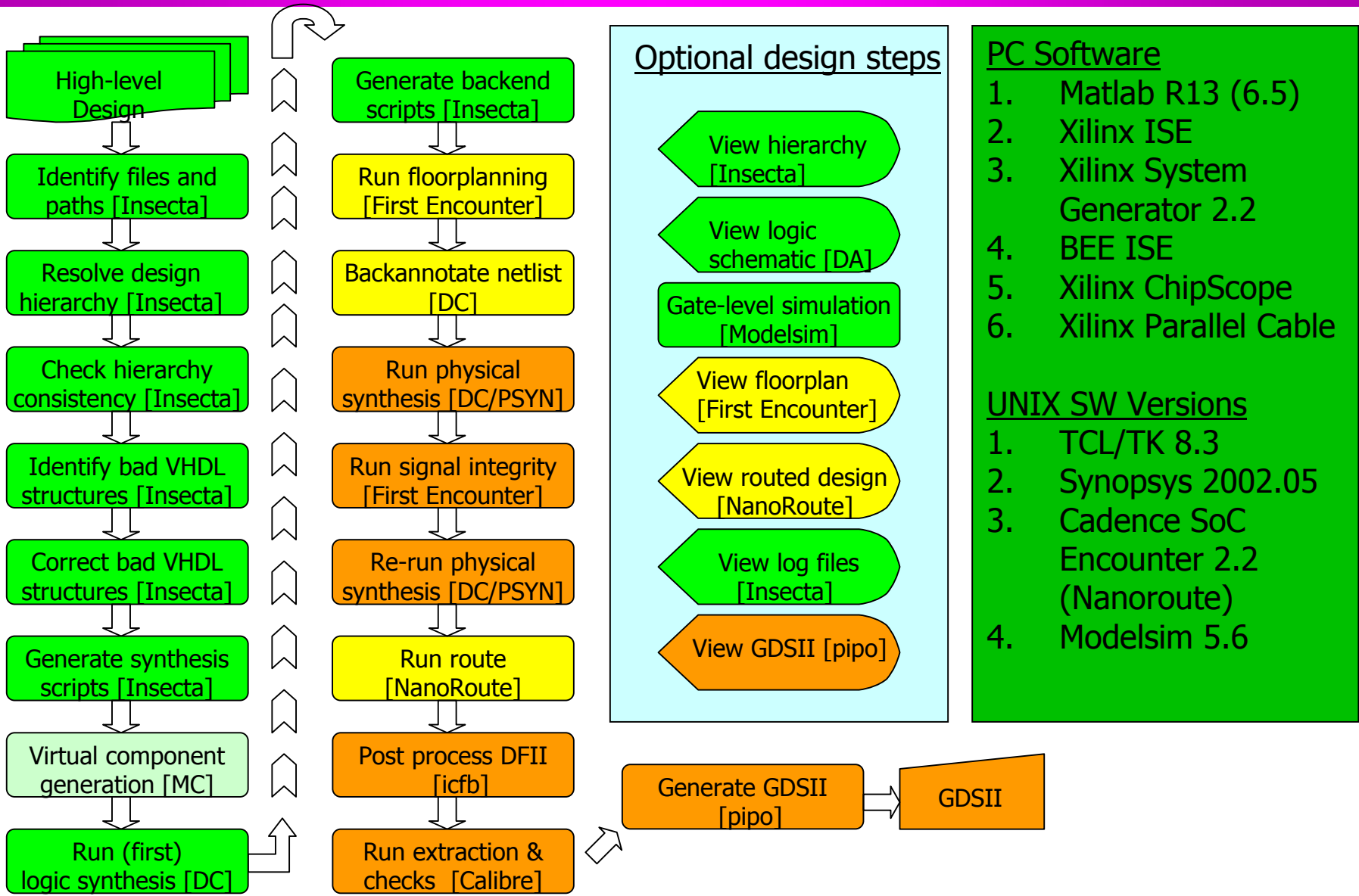
- Tcl/Tk code drives the flow
  - » Used to drive multiple EDA tools: First Encounter, Nanoroute, Module Compiler
- GUI controls technology selection, parameter selection, flow sequencing
  - » A real “Push Button” flow...
  - » Users can refine flow-generated scripts

The screenshot shows the INSECTA GUI with the following configuration:

- Working directory: /vol/hitz/vol2/designs/sshaf/users/richards/demo/SYSGEN/
- Top-level filename: /tools/sshaf/users/richards/demo/SYSGEN/demo\_subsystem.vhd
- Effort: Low
- Optimize for 100MHz clock
- Start: run\_first\_place
- Tech: ST 0.13u LL Worst
- Stop: end
- Hierarchy: Boundary Optimize
- Advanced Flow...
- Target: FirstEncounter
- Status: idle

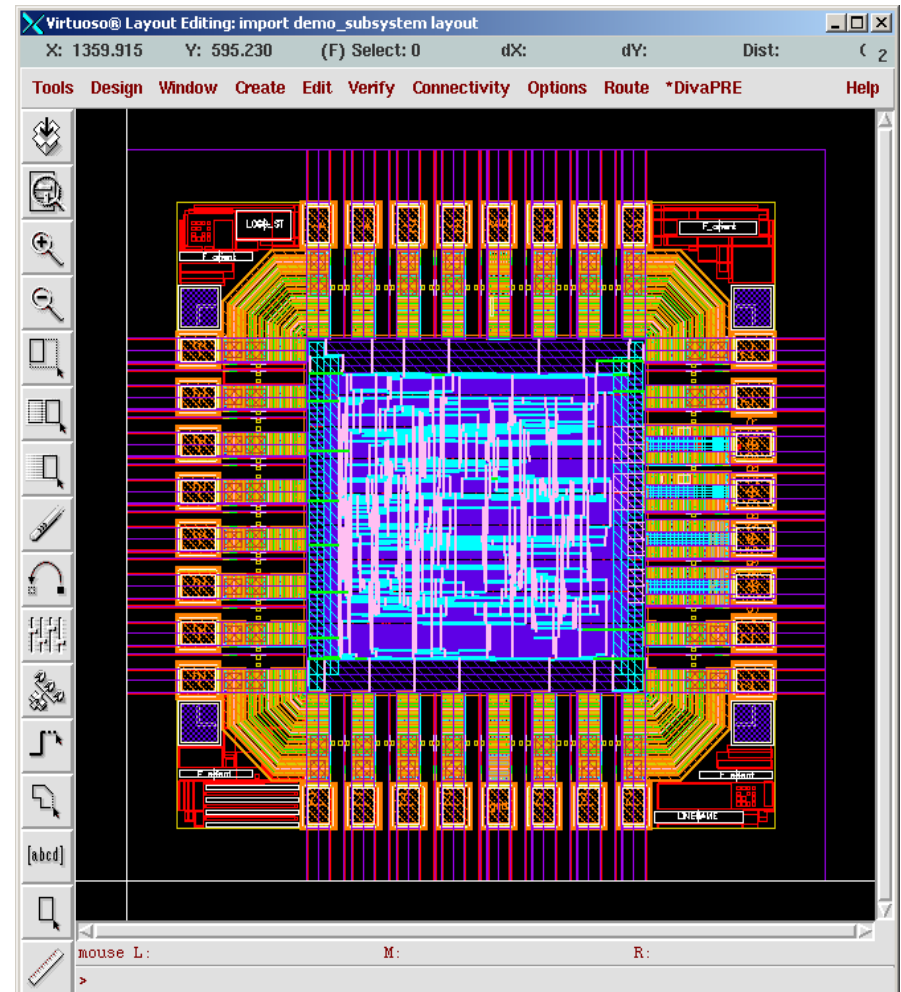
Buttons: Debug, Help, Exit, Run INSECTA

# Automated ASIC flow tools



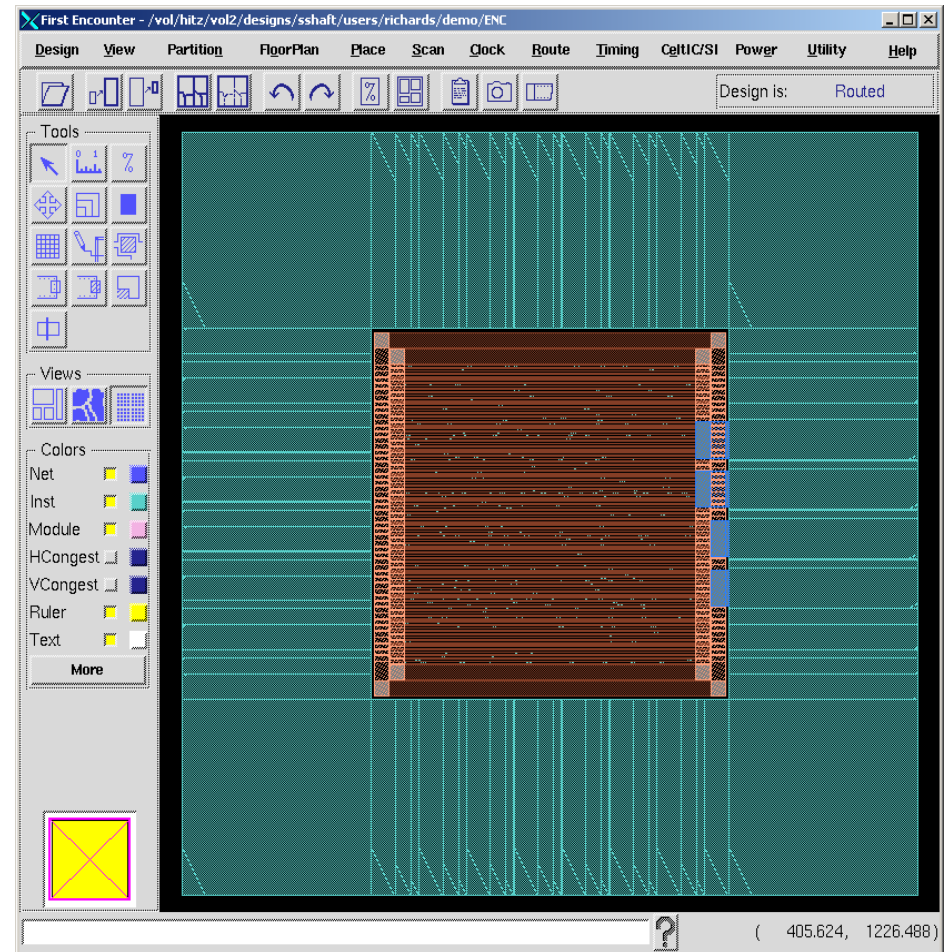
# ASIC Flow: Back-end

- Using Unicad (ST Microelectronics) backend directly for DRC, LVS, Antenna rule checking
  - » Easier to track technology updates from foundry.
  - » Critical for evaluating internally developed technology files for FE, Nanoroute



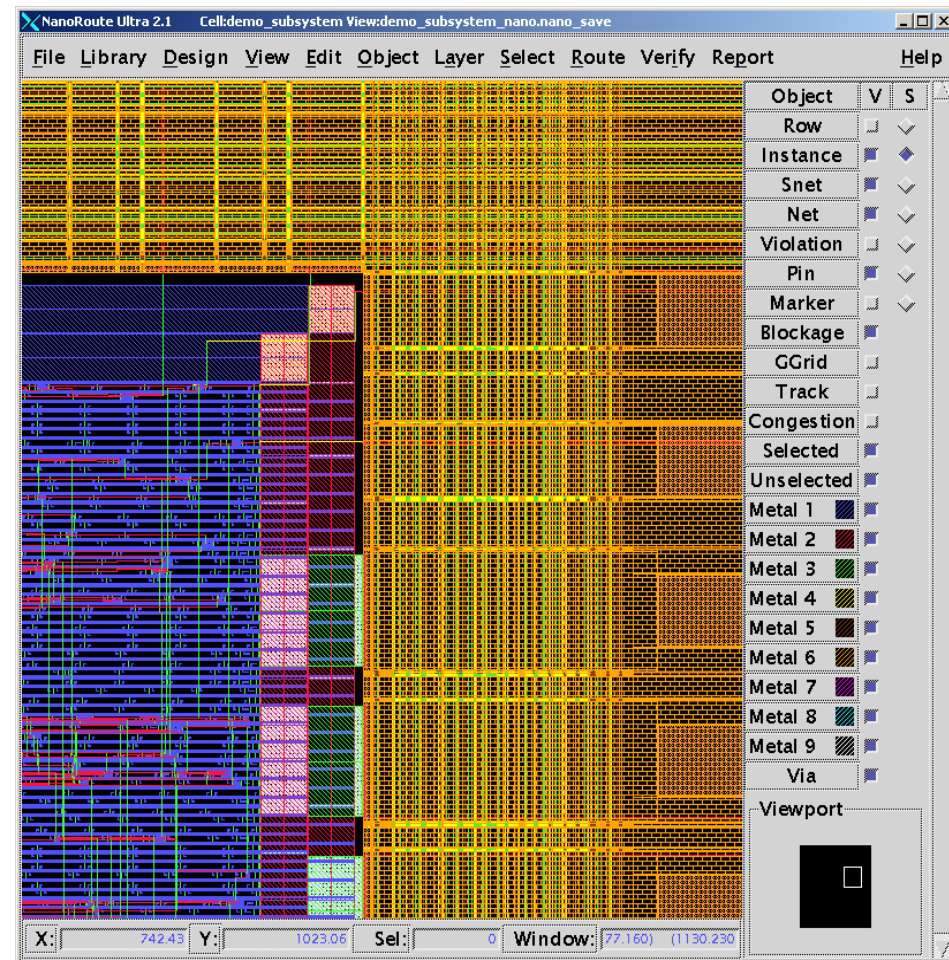
# ASIC Tool Flow: Placement

- Cadence based flow
  - » First Encounter (FE)
  - » Nanoroute
- Timing Driven!
  - » FE provides accurate wire parasitic estimates
  - » Placement by FE



# ASIC Flow: Routing in 130nm

- Nanoroute: Ready for 130nm, 90nm designs
  - » Stepped metal pitches
  - » Minimum area rules
  - » Complex VIA rules
  - » Avoids antenna rule violations
  - » Cross-talk avoidance: to be evaluated
- Silicon Ensemble: Fallback position
- Apollo tools: Possible alternative



# ASIC directly from Simulink – Narrowband Transmitter

CPU time: 57 min

Core Utilization: 0.344418 (Pad limited)

Size (From SoC Encounter):

Core Height 565.8u

Core Width 489.54u

Die Height 1322.66u

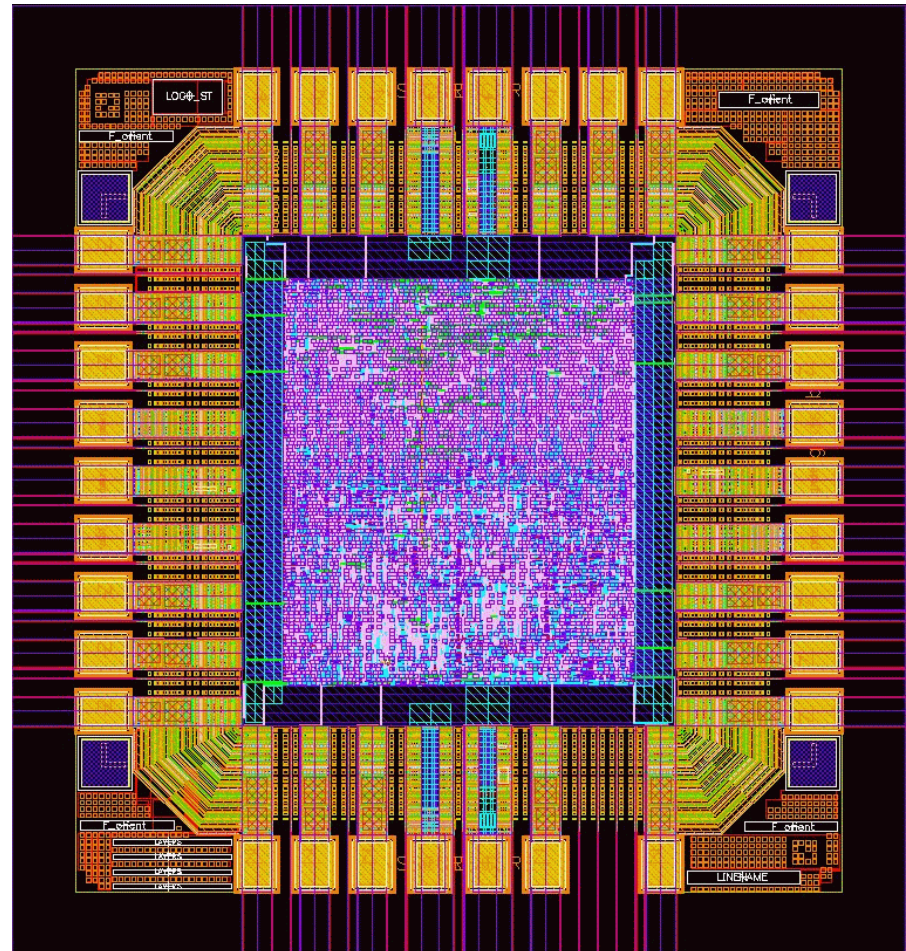
Die Width 1242.3u

Synopsys estimates:

Total Dynamic Power = 610.5163  
uW (100%)

Cell Leakage Power = 15.9364 uW

Critical path: 9.21ns



# The Issues I Addressed

---

- How much flexibility is needed and how best to include it...
  - » As little as possible consistent with business constraints
- A single system description including interaction between the analog and digital domains
  - » Timed dataflow plus state machines
- “Realtime” SOC prototyping
  - » FPGA configurability makes real-time prototyping possible in a fully parallel architecture.
- Automated ASIC design flow
  - » Certainly possible - the “chip in a day” flow