

MTAAP'07 Keynote

Michael Merrill



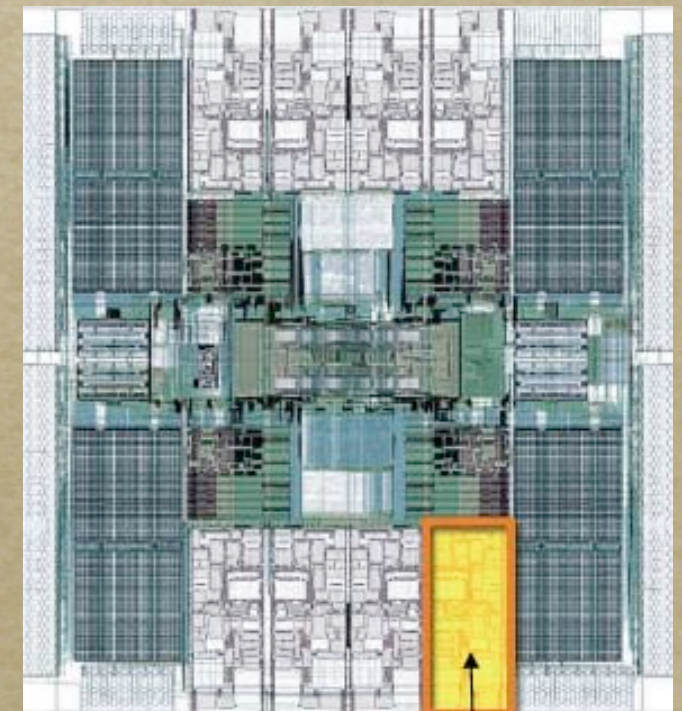
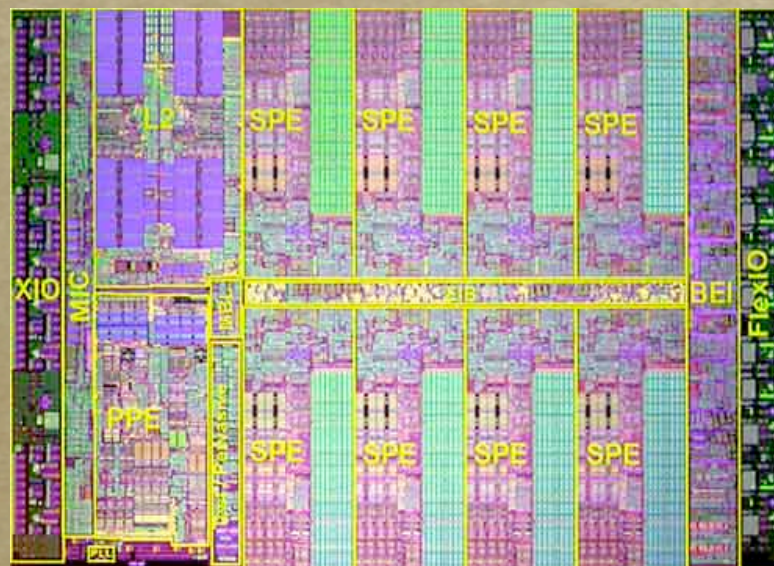
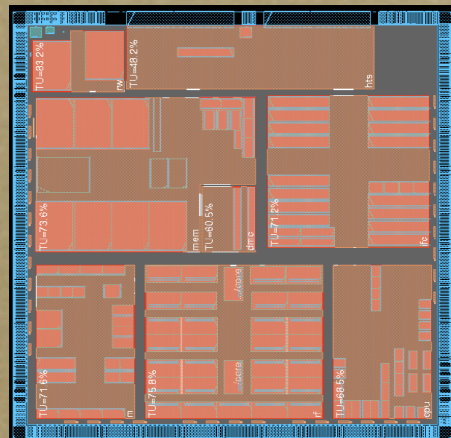
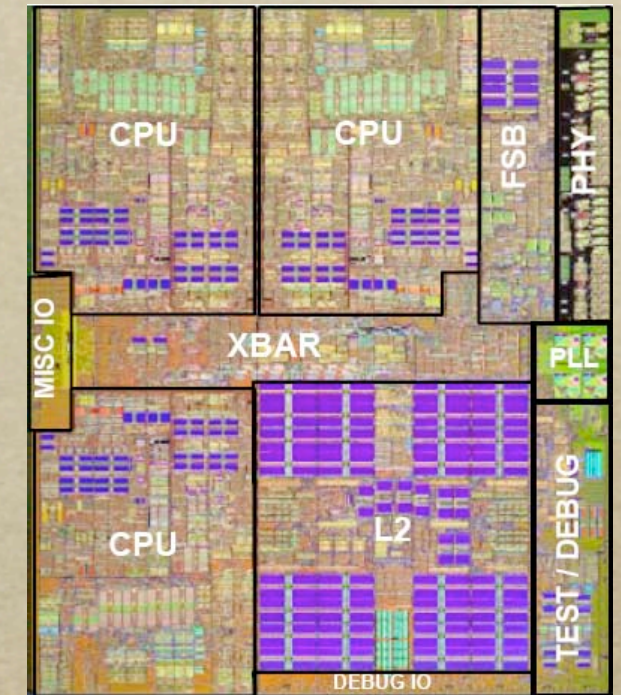
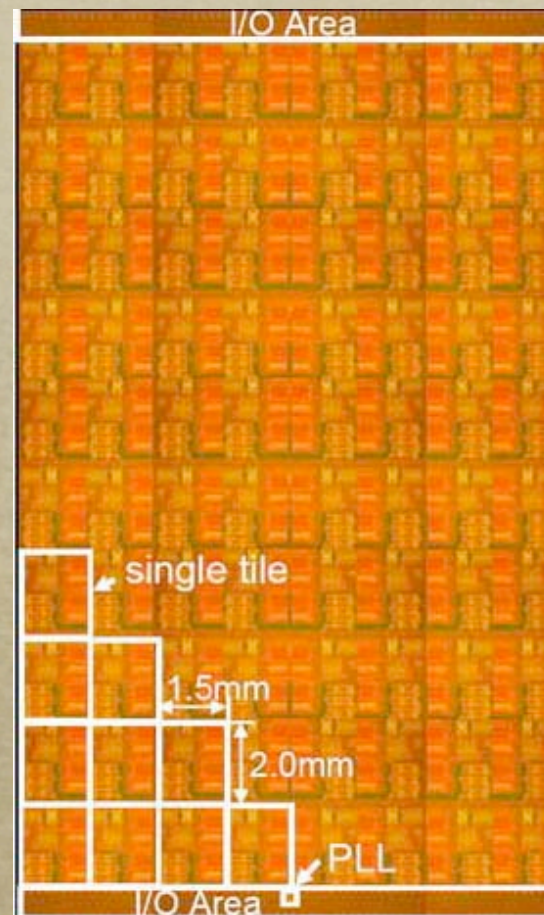
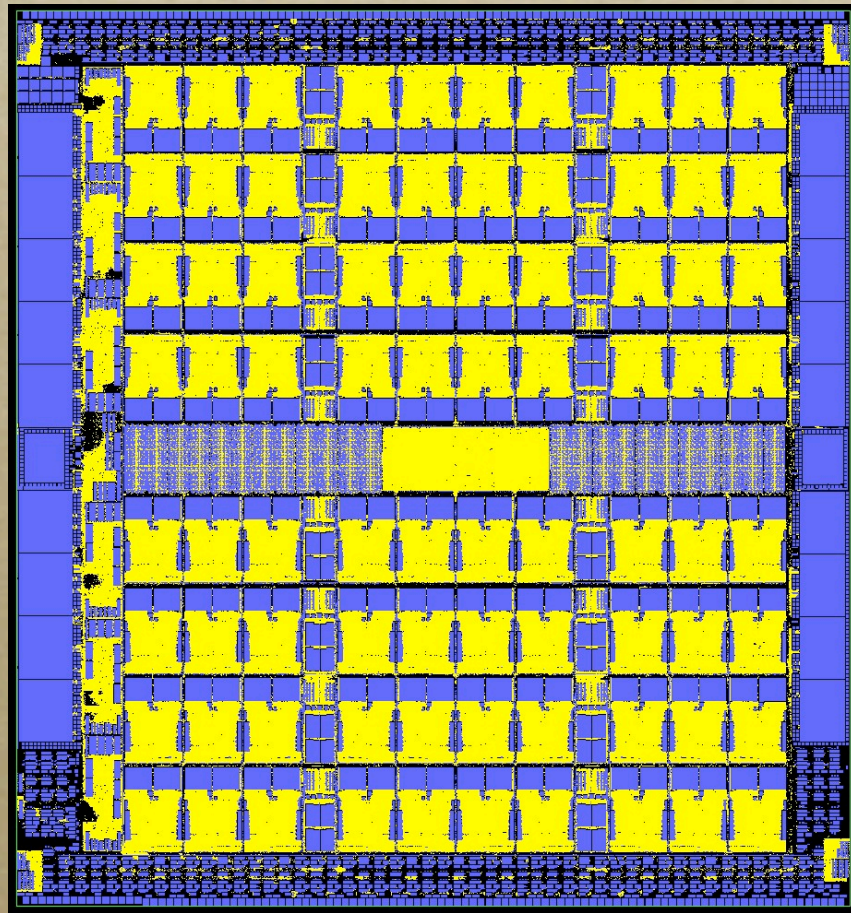
Outline

- *What's important... applications*
- *Making sense... of all this stuff*
- *What's necessary... I think*
- *What's possible... maybe*

What's important... applications

- *Data Structures are important... more support for linked data structures*
- *Ignored algorithm areas are coming back to bite us*
- *Sparse methods on unstructured data*
- *Adaptive methods are better aligned with nature but not with current architecture*
- *Helping humans deal with information overload*

Making sense... of all this stuff



Spectrum

*Hardware contexts per
set of functional units*

*all contexts to
one set of
functional units*

*one context to
one set of
functional units*

MTA/XMT

UltraSparc T1

Cyclops64

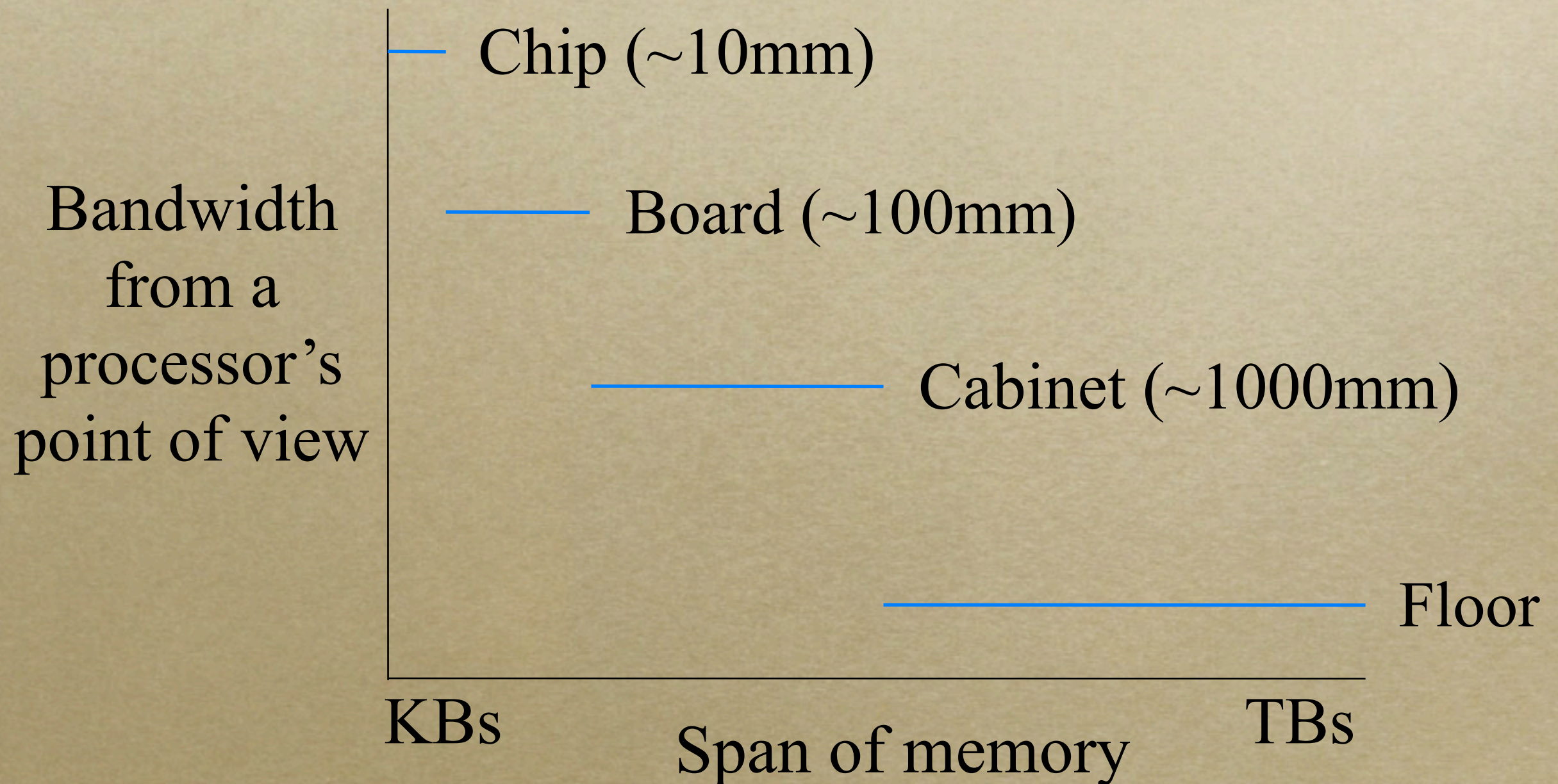
Stuff

- *Virtualization*
 - *How much is enough?*
 - *Fault tolerance*
- *How much baggage does a context have?*
Probably affects virtualization
- *Synchronization*

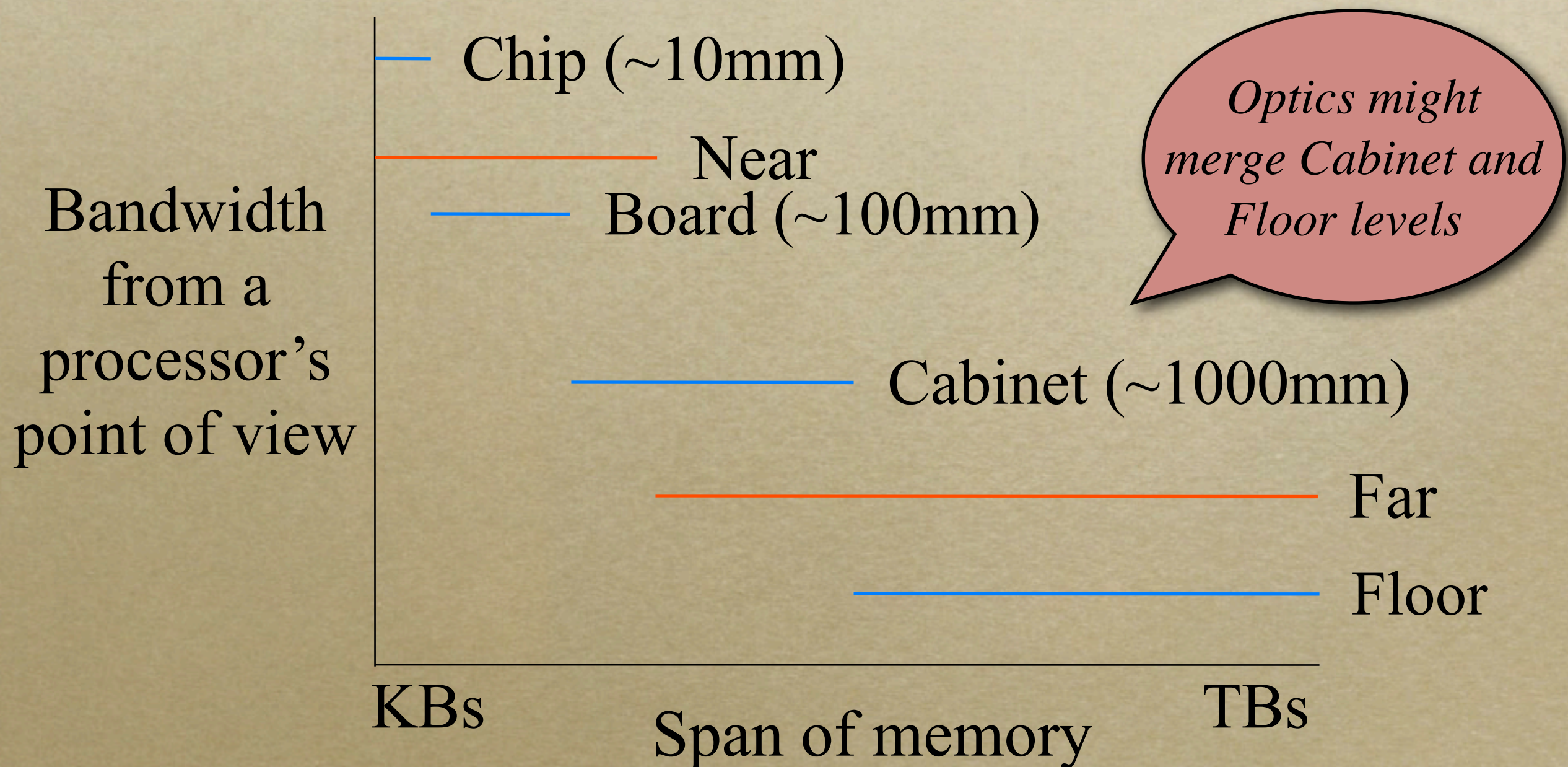
More Stuff

- *Explicit memory hierarchy?*
- *I-cache! Don't make the programmer worry about code size!?!*
- *Commercial use vs. scientific use... vs. something else*

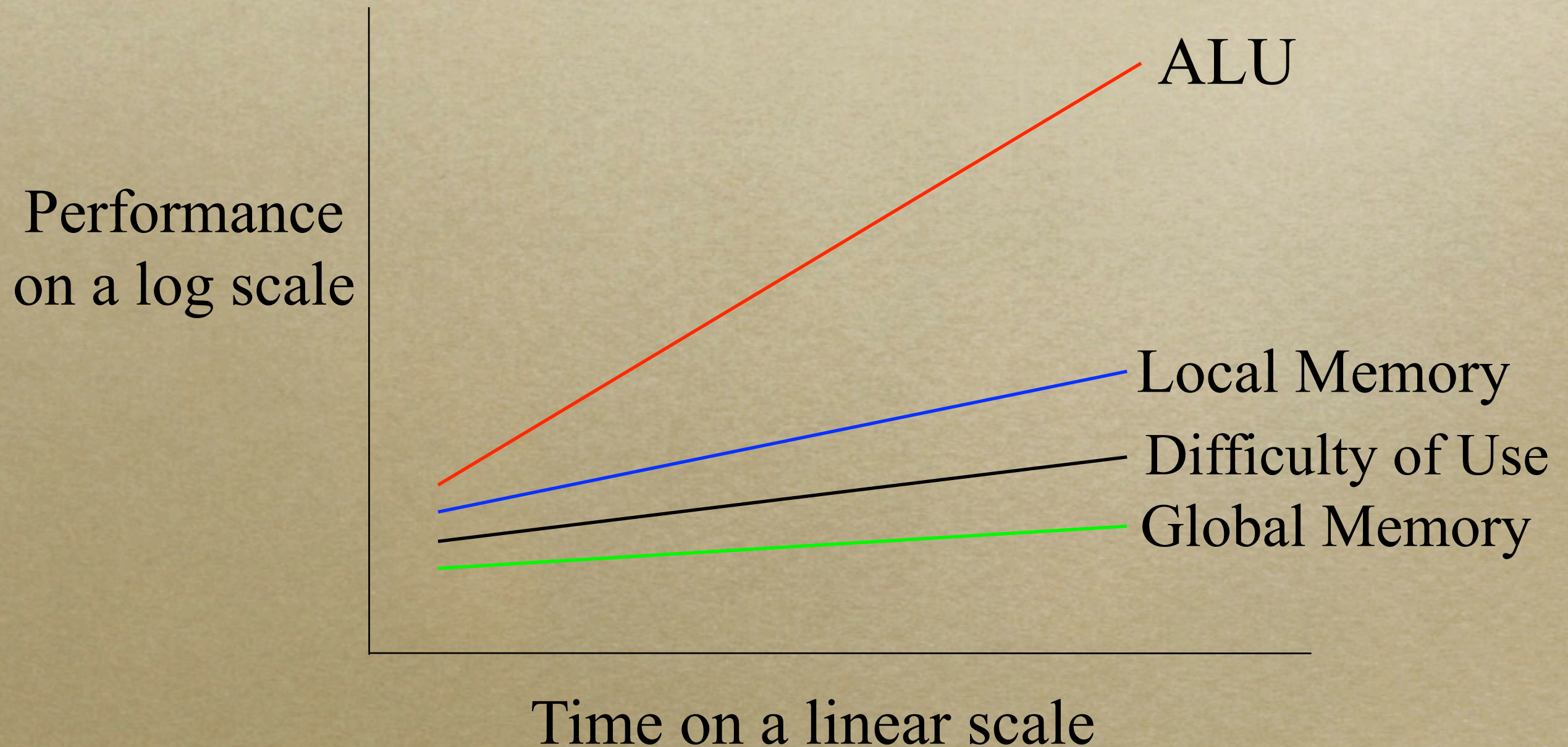
Natural Bandwidth Boundaries



Natural Bandwidth Boundaries



Trends we live with...



Different Balance

- *Costs have changed drastically*
- *Transistors are cheap... Wires are expensive*
- *Processor complexity vs. power is an issue*
- *Balance costs... apply transistors to use wires more effectively... not just for cache*
- *This is why you see architecture changing*

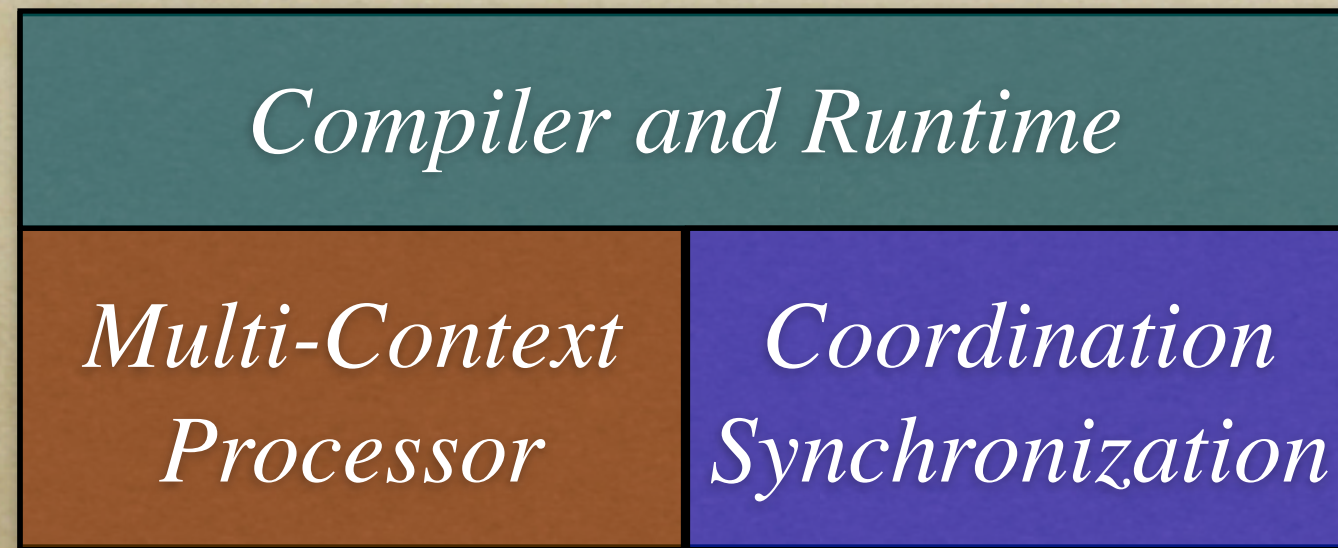
What's necessary... I think

- *Need to provide an effective system solution HW and SW!*
- *Why? Days of coarse grained scaling are at an end... so threads/contexts will necessarily work together to perform a task.*

Fabrication

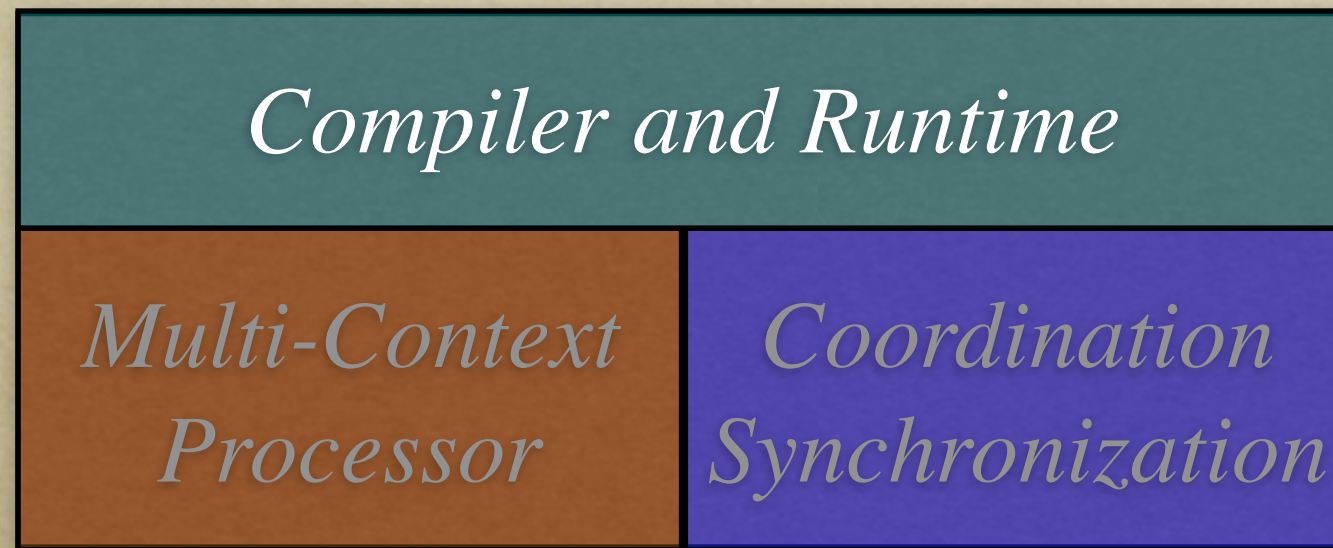
- *Tiled architecture with partial good chips for lower costs*
- *Detect failed computation*
- *Retry failed computation*
- *Move away from fixed number of threads/ contexts*

What's necessary... I think



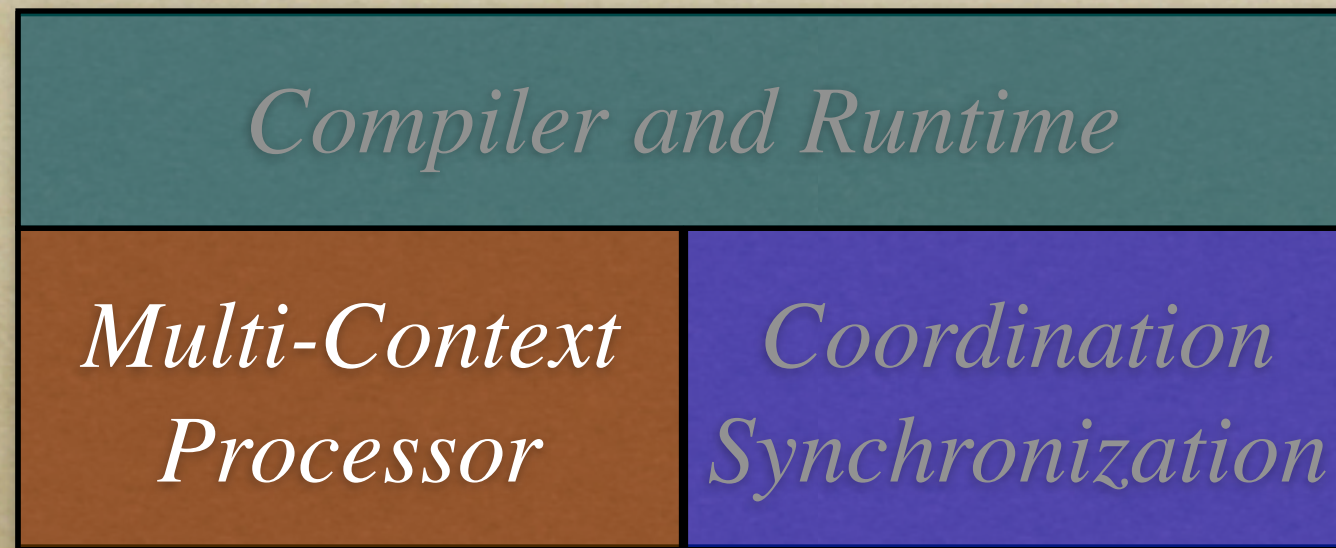
- *Need at least these three things working together to produce an effective environment for the application developer*

Effective Software



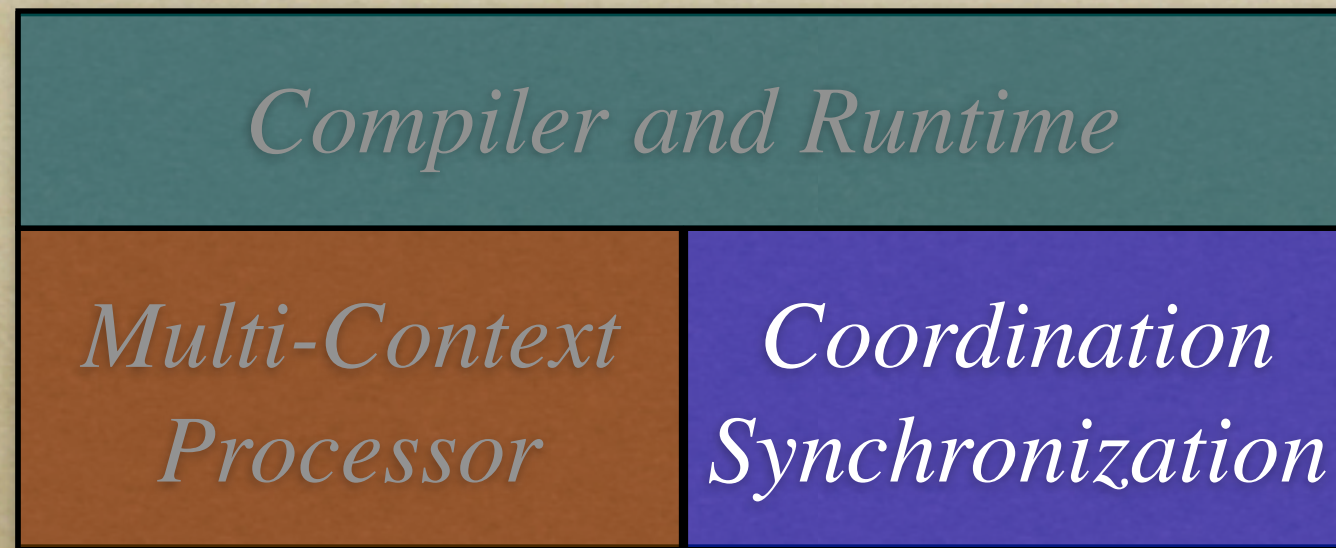
- *Runtime that provides effective dynamic work management so the unbalanced nature of the application can be mitigated.*
- *Compiler that takes advantage of such a runtime increases programmer effectiveness and productivity allowing them to concentrate on the application.*

Latency Tolerance/Management



- *Effective use of the bandwidth provided by the internal system networks through the use of latency tolerance and/or latency management techniques.*
- *Many of these techniques require the exposure of abundant fine-grained parallelism in the application.*

Low Overhead Coordination



- *Threads will necessarily work together to compute so effective coordination will be essential.*
- *Any cycles spent waiting on synchronization events are not spent computing and therefore decrease efficiency.*

What's possible... maybe

- *Don't look for any major companies to make things significantly better because it messes with the current business too much.*
- *Which direction to go?*

Straight Forward Scaling

	<i>90nm</i>	<i>65nm</i>	<i>45nm</i>	<i>32nm</i>
<i>TU</i>	<i>160</i>	<i>306</i>	<i>640</i>	<i>1266</i>
<i>FPU</i>	<i>80</i>	<i>153</i>	<i>320</i>	<i>633</i>
<i>TU/XB</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>6</i>
<i>XBar</i>	<i>80</i>	<i>102</i>	<i>160</i>	<i>211</i>
<i>Clock</i>	<i>500M</i>	<i>585M</i>	<i>684M</i>	<i>800M</i>
<i>Perf</i>	<i>80G</i>	<i>179G</i>	<i>437G</i>	<i>1.01T</i>
<i>SRAM</i>	<i>4.8M</i>	<i>9.2M</i>	<i>19.2M</i>	<i>37.9M</i>

- *start with Cyclops64*
- *22x23mm die*
- *150W to 190W*
- *3DE ?*

What About Software?

- *Need good compiler technology to exploit on chip explicit memory*
- *Much higher level of abstraction*
- *Need to separate the how-to from the what-for but express both*
- *Diagnose hot spots (resource contention)*
- *etc...*

- *Questions? ... I have a ton ;-)*