

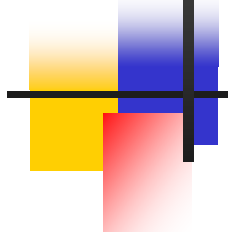


Hybrid Computer Architecture

Brian Van Essen

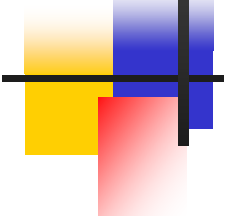
Benjamin Ylvisaker

Carl Ebeling



Moore's Law: Is it Over?

- n von Neumann processors no longer scale
 - n Overhead of speculative execution is too high
 - n Complexity of superscalar OOO core is n^2
 - n Optimum power / performance pipeline depth is ~ 7 stages
- n Spatial processors benefit from added transistors
 - n Reconfigurability allows virtualization
 - n Enables programming abstraction



Keeping up with streams is hard

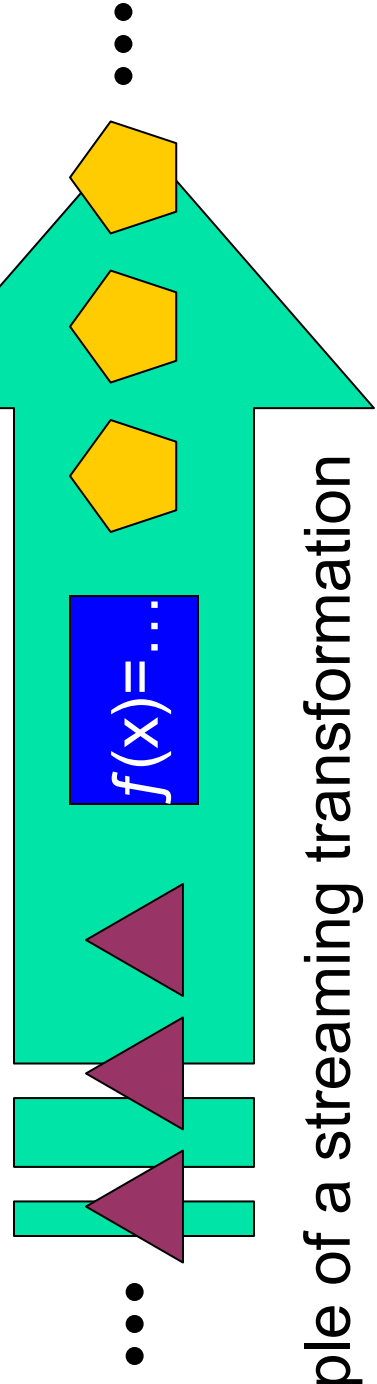
n Multimedia workloads

n Audio & Video

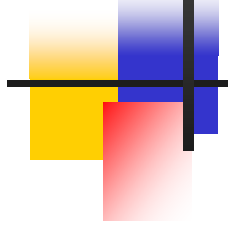
n Communication workloads

n Networking

Spatial processors
are good at this

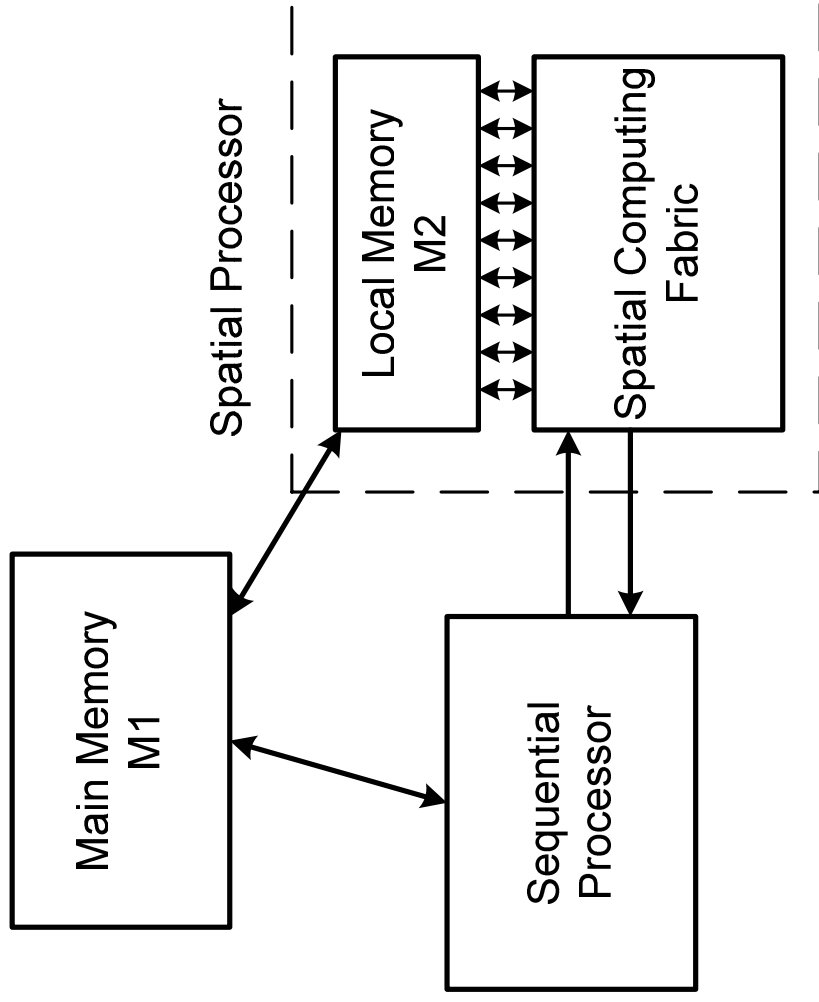


Example of a streaming transformation



Hybrid Architecture Research

- Blend sequential and spatial computing
- One program executes both types of computation





Overview

- n What is spatial computing
 - n Why is it interesting
- n Hybrid Architectures
 - n What is hard about hybrid architectures
- n Future Research



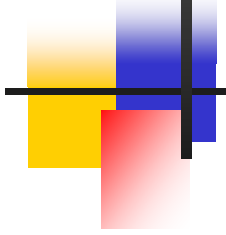
What is spatial computing?

n Spatial processors:

- n** Parallel array of compute elements (fabric)
- n** Assign operations to different physical resources
- n** Stream operands through the fabric
- n** Execute many operations in parallel

n Sequential processors:

- n** Step through a sequence of instructions

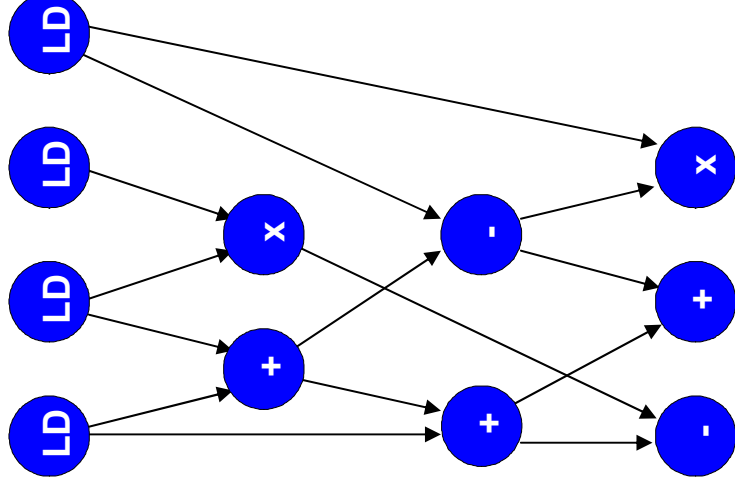


Encoding a program

Instruction Stream

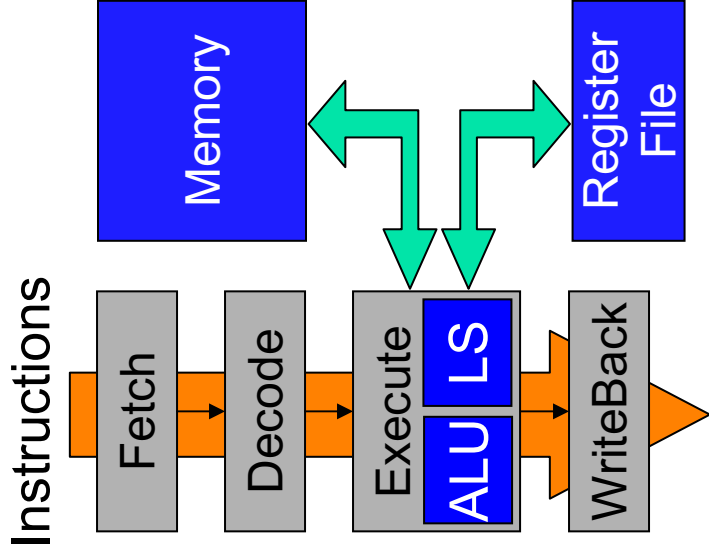
Load r1, A
Load r2, B
Load r3, C
Load r4, D
Add r5, r1, r2
Mul r6, r2, r3
Add r7, r1, r5
Sub r8, r5, r4
Sub r9, r7, r6
Add r10, r7, r8
Mul r11, r8, r4

Dataflow Graph

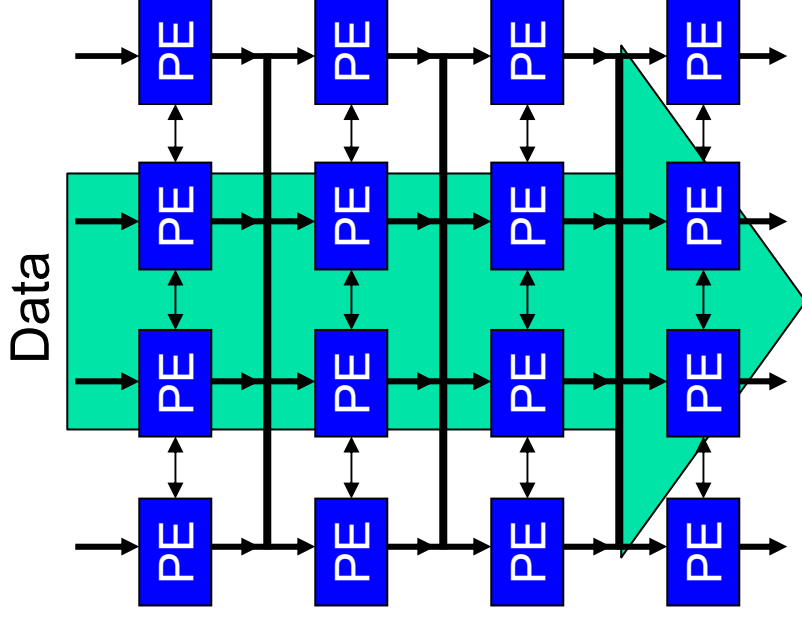


Processors: Under the hood

Traditional Computer
(Load / Store Arch)



Spatial Computer
(e.g. FPGA, PipeRench)





Why spatial processors?

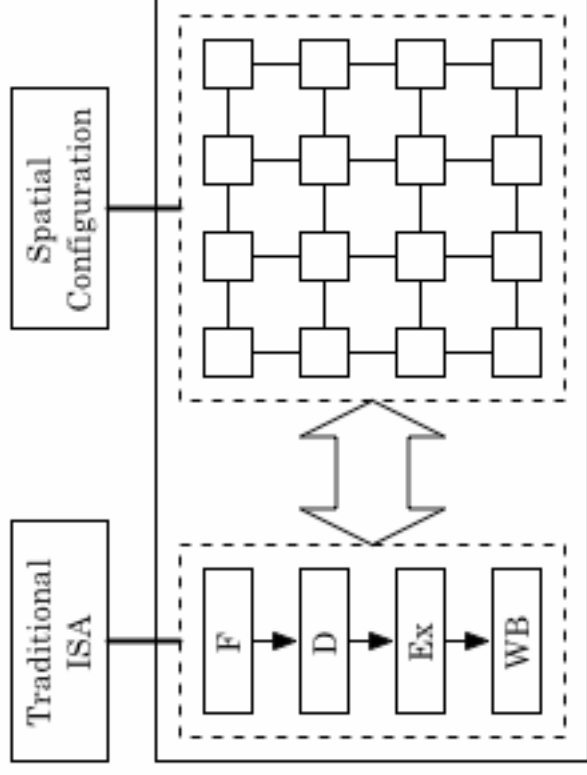
- n Extremely efficient for certain applications
 - n Regular computation
 - n Regular communication
 - n e.g. Streaming Data
- n Excellent performance / power ratio
- n Limitations:
 - n Difficult to execute control flow
 - n Hard to program

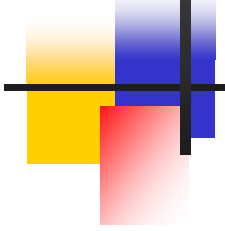
Basic Hybrid Architectures

- Two processors on a single chip
 - Integrates control plane and data plane processors
 - Provide high speed interconnect
 - Share memory

Execute independent programs

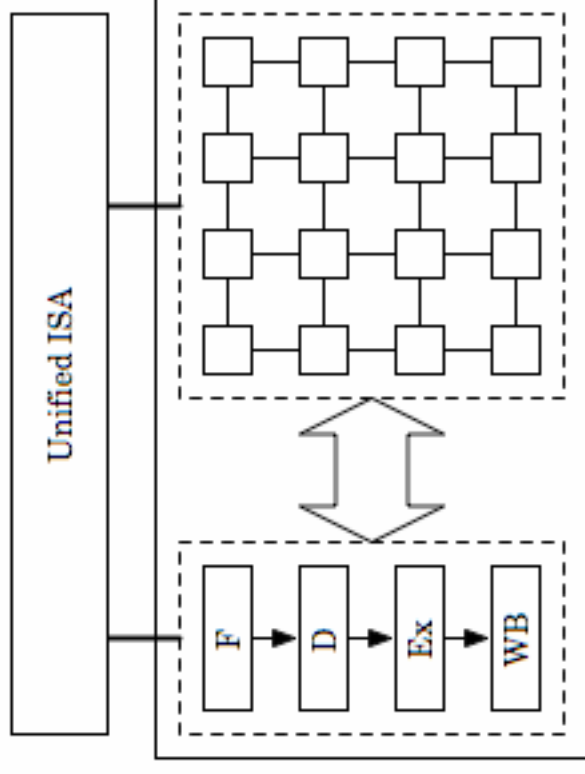
- Manage synchronization

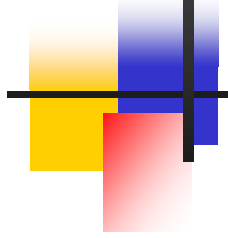




Unified Hybrid Architecture

- n Single programming model
 - n Collapses control plane and data plane processors into single abstraction
 - n Implicit synchronization
 - n Simplified programming abstraction
 - n Program “Automagically” executes on appropriate processor
 - n Runtime system manages fabric configuration





Research Challenges

Creating a new Instruction Set Architecture (ISA)

- Provides canonical sequential interpretation
- Exposes good spatial configuration
- Efficient synchronization of runtime control

Virtualization of spatial processors is hard

- Necessary to provide abstract programmers model
- Use dynamic reconfiguration

Programming Language

- Explicit stream operations
- Disambiguate memory references



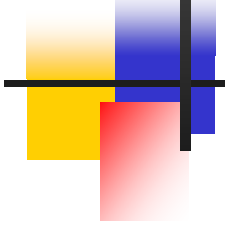
Research Synopsis

- n Define new processor architecture and ISA
 - n New level of ease of use
 - n Unified programming model
 - n Blend sequential and spatial computing
 - n Excels at streaming data applications
 - n One program executes both types of computation
 - n Implicit communication
- n Efficient virtualization of spatial processors
- n System-level programming language



Appendix

Type Architectures Programming Languages

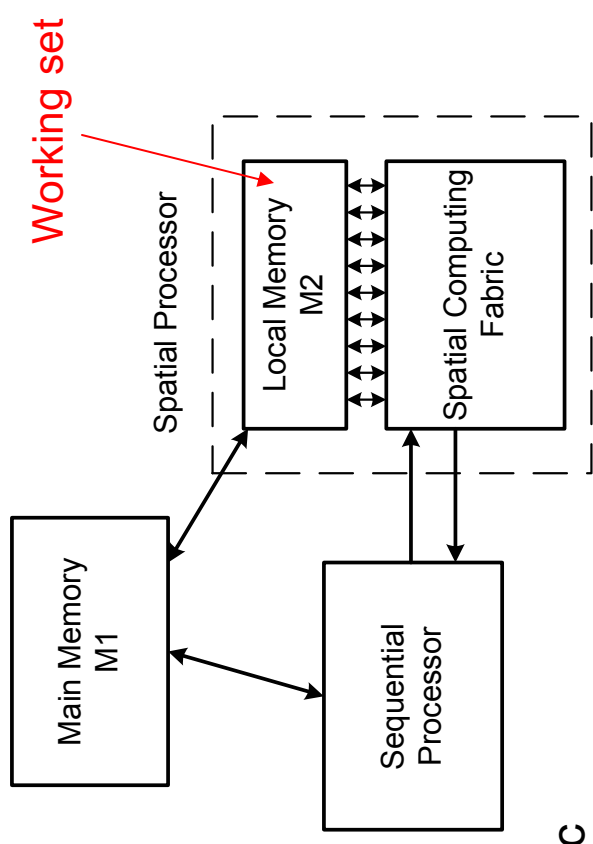


Abstract processor models

- n von Neumann Type Architecture - RAM Model
 - n A processor interpreting 3-address instructions
 - n PC describing the next instruction of program in memory
 - n Flat, randomly accessed memory requires 1 time unit
 - n Memory is composed of fixed sized addressable units
 - n One instruction executes at a time, and is completed before the next instruction executes
- n Modern RISC & CISC processors emulate this model
 - n C directly implements this model

Hybrid Type Architecture

- von Neumann sequential processor
 - Spatial Fabric
 - P operations per cycle
 - Statically scheduled
 - Main Memory
 - ~ 1 access per cycle
 - Local Memory (Workspace)
 - ~ P accesses per cycle
 - enough to maintain P ops
 - Alternating Execution
 - Sequential program executes
 - Control transferred to spatial fabric
 - Shared state transferred
 - Atomic execution of spatial section
 - Shared state transferred back





A new Programming Language

- ┆ “System level”

- ┆ Full control of underlying ISA
- ┆ Explicit resource management

- ┆ Key Issues

- ┆ Expressing parallel portions of computation
 - ┆ Easily mapped to spatial processor
- ┆ “Relaxed” memory access ordering
 - ┆ e.g. streams
- ┆ Disambiguate memory references
 - ┆ mitigate aliasing
- ┆ Reflect constraints of type architecture
 - ┆ e.g. low main memory bandwidth