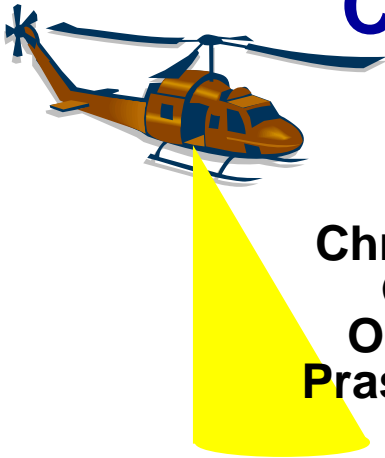




CEARCH

Cognition Enabled ARCHitecture



Stephen Crago and Janice McMahon, USC/ISI

Chris Archer¹, Krste Asanovic², Richard Chaung³, Keith Goolsbey⁴, Mary Hall⁵, Christos Kozyrakis⁶, Kunle Olukotun⁶, Una-May O'Reilly², Rick Pancoast⁷, Viktor Prasanna⁸, Rodric Rabbah², Steve Ward², Donald Yeung⁹

September 20, 2006

¹Northrop Grumman, ²MIT, ³Army I2WD, ⁴Cycorp, ⁵USC/ISI, ⁶Stanford University, ⁷Lockheed Martin, ⁸USC, ⁹University of Maryland

The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the Defense Advanced Research Projects Agency (DARPA) or the U.S. Government. Effort sponsored by the Defense Advanced Research Projects Agency (DARPA) through the Department of the Interior National Business Center under grant number NBCH104009.



Outline



- **Project Goals**
- **Architecture Characteristics**
- **Application Examples**
- **Summary**



CEARCH Goals



- **Develop a *computer architecture* that supports cognitive information processing**
 - Computer architecture: a set of hardware and system software interfaces and implementations
- **Support real-time, embedded cognitive processing requirements through an efficient, high-performance computer architecture**
- **Identify algorithms and improved algorithm implementations that can leverage the CEARCH computer architecture**

- **CEARCH is not a *cognitive architecture* project**
 - Cognitive architecture: a computational model (usually expressed in software) for a complete cognitive system that may or may not be based on human psychology



CEARCH and Cognitive Architectures



- **The CEARCH computer architecture will run a *variety* of cognitive architectures efficiently**
 - **Multiple cognitive architectures important**
 - No single consensus on cognitive architectures
 - Important to support emerging cognitive architecture research: each IPTO program in this domain has its own cognitive architecture
 - Different domains may require different cognitive architectures
 - **Support for variety of cognitive architectures**
 - Wide range of cognitive algorithms drive CEARCH architecture to ensure coverage
 - Adaptivity and scalability emphasized to support dynamic processing requirements critical to all cognitive architectures

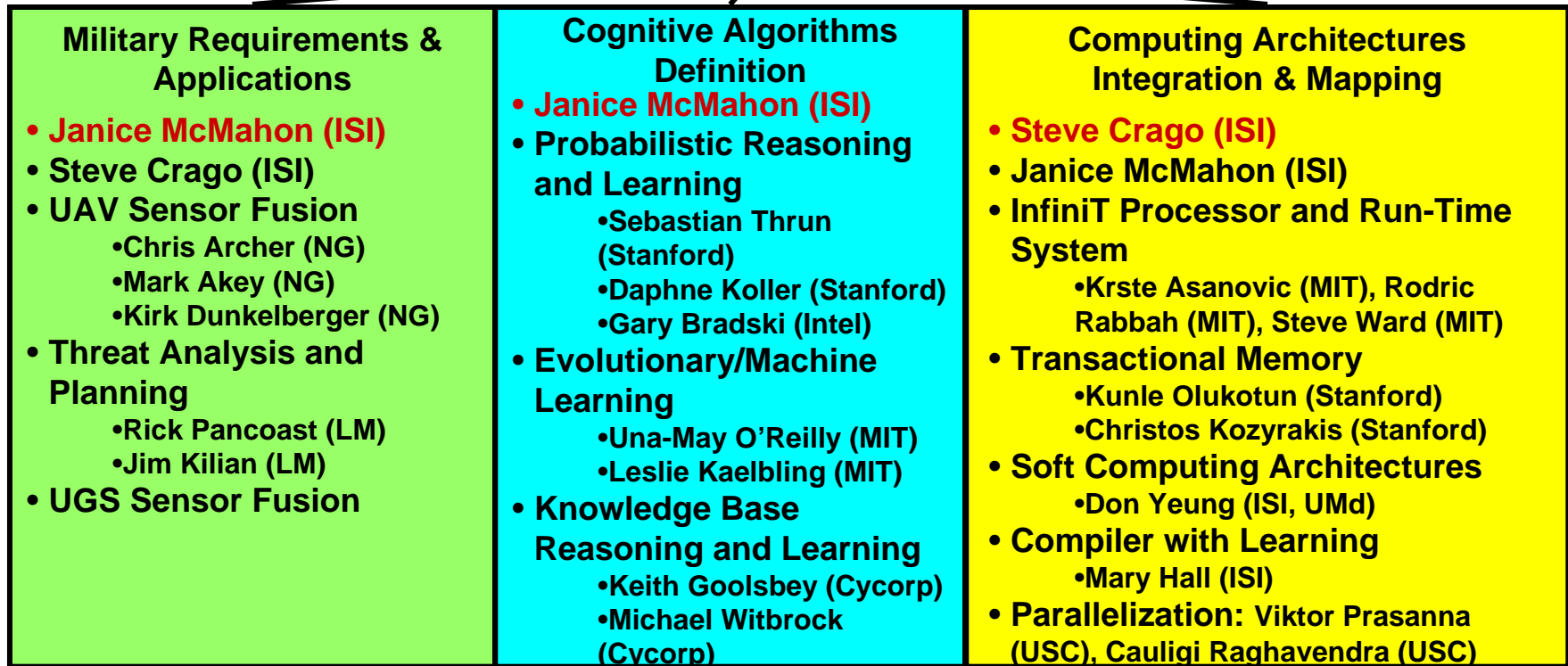
- **CEARCH computer architecture has some characteristics of a cognitive system**
 - **Introspection and self-management: knows what it is doing and how to process efficiently**
 - **Learns how to process more efficiently over time**
 - **Supports inexact computations when optimality is not feasible or possible**
 - **Robust processing in the context of faults**



CEARCH Team



Program Lead
 Steve Crago (Co-PI, ISI)
 Janice McMahon (Co-PI, ISI)
 Bob Parker (ISI)



STANFORD UNIVERSITY



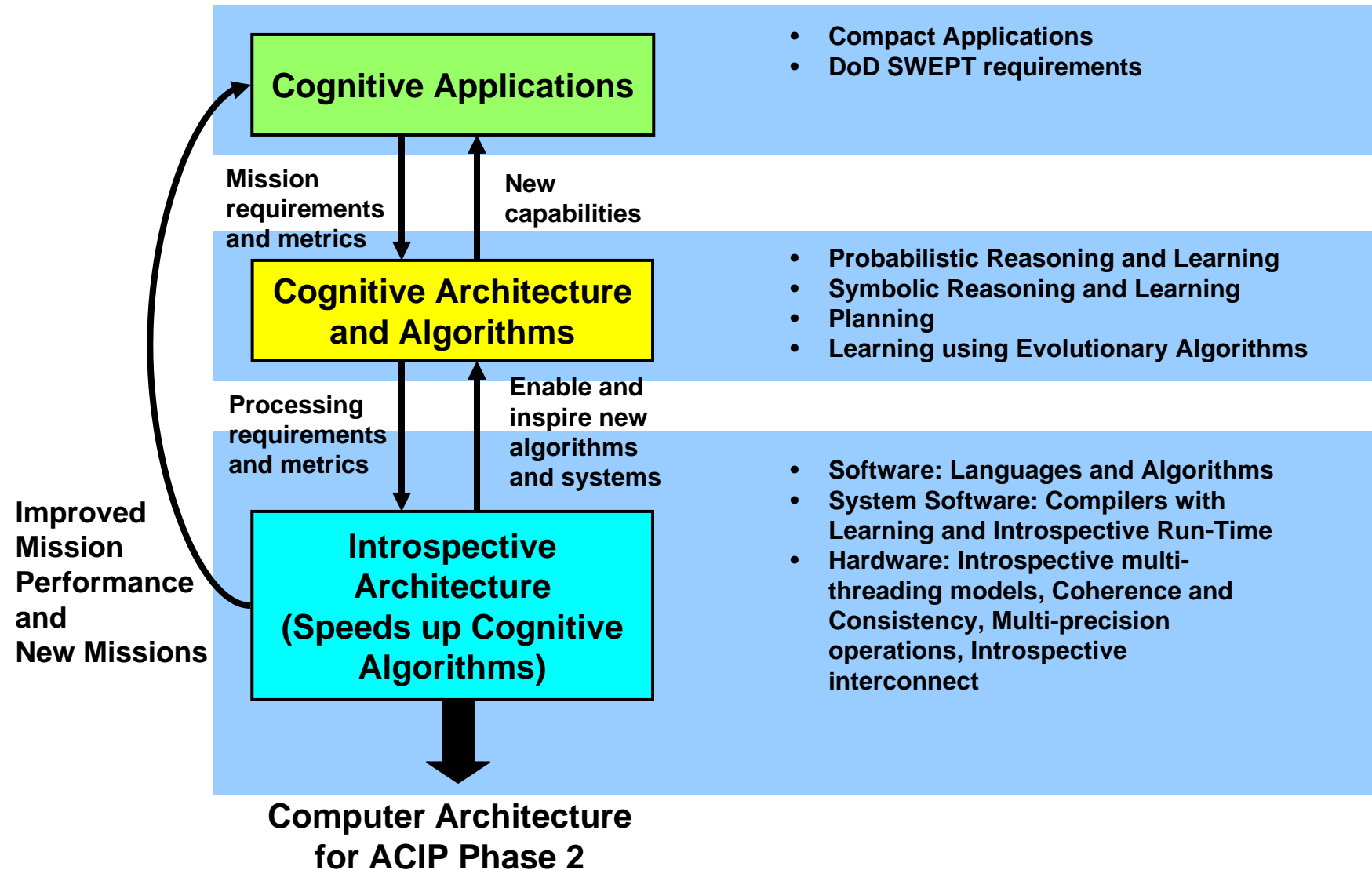
USC Viterbi School of Engineering

STANFORD UNIVERSITY





CEARCH Project Overview

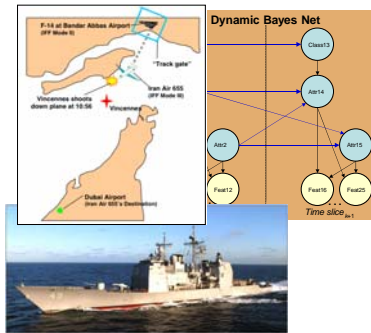




Scenario Summary



Shipboard Threat Analysis and Planning



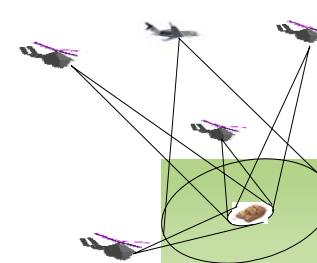
UGS Urban Situational Awareness



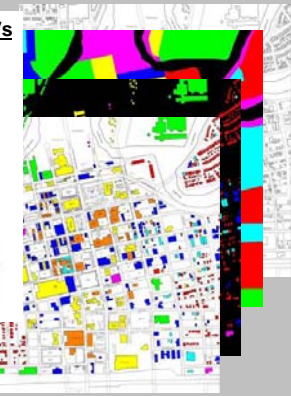
LW-451

UAV-based Behavior Spotting

Multi-UAV Sense/Attack Scenario



Autonomous UAVs



Kernel	Example Scenario Requirement	Example architectural drivers
Probabilistic Relational Model (Learn, Infer)	1-2 Tera-updates / sec on large graphs	Probabilistic computation
SATisfiability-based Planner	1 Giga-Boolean-inferences / sec	Parallel tree traversal
Support Vector Machine Classification	2 Tera-ops (variable-precision floating point) / sec	Flexible caching for sparse vectors
Information-form Data Association Tracking	2 Tera-ops (probability calculations) / sec	Parallel sparse matrix calculations
Symbolic Reasoning and Learning	313K problem trees per second	Symbolic matching, irregular memory accesses
System		Rapid High-Level Reorganization and Responsivity

Cognitive reasoning and learning techniques require new computing platforms to enable new real-time, embedded capabilities and missions
Must combine orders of magnitude performance/efficiency improvement with ability to respond rapidly to the needs of dynamic environments



Outline



- **Project Goals**
- **Architecture Characteristics**
- **Application Examples**
- **Summary**



Why Do We Need Hardware for Cognitive Systems?



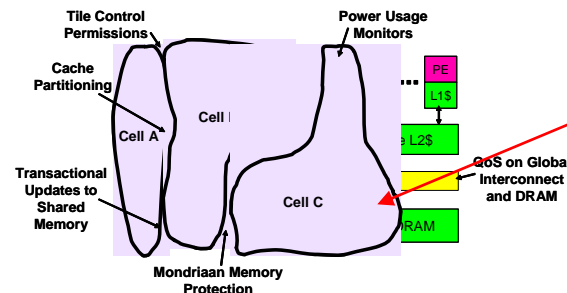
- **Introspective and Self-Managing Computing**
 - Must support introspective information flow from applications to hardware (and back) to support cognitive resource management and introspective applications
 - **Scalable Web of Cognitive Virtual Processing Elements**
 - Efficient, high-performance computation required to support real-time reasoning and learning requirements
 - Must be adaptable and able to support variety of cognitive processing paradigms (graphs, symbolic reasoning, etc.) and dynamic requirements
 - **Multi-level Soft Computing**
 - Support for probabilistic and inexact data types and computation pervasive in system (processing, memory, communication, programming model, run-time system)
 - **Adaptive memory system**
 - Unpredictable, irregular memory accesses and large working sets
 - Driven by parallel computation, dynamic resource allocation, and fundamental characteristics of algorithms and data



Introspection and Self-Management



- **System must adapt to unpredictability in cognitive systems**
 - **Dynamic scenarios lead to dynamic and unpredictable changes in processing requirements**
 - **Cognitive processing too complex to be managed by programmer**
 - Cognitive algorithms provide means for system to manage itself
 - **Faults are unavoidable at this scale**
- **Introspection required to support autonomous adaptability**
 - **Processing:** precision, performance required, operation mixes, efficiency of functional units
 - **Memory and Communication:** access/communication patterns, cache hit rates, working set sizes, precision required, bandwidth/latency trade-offs, protection



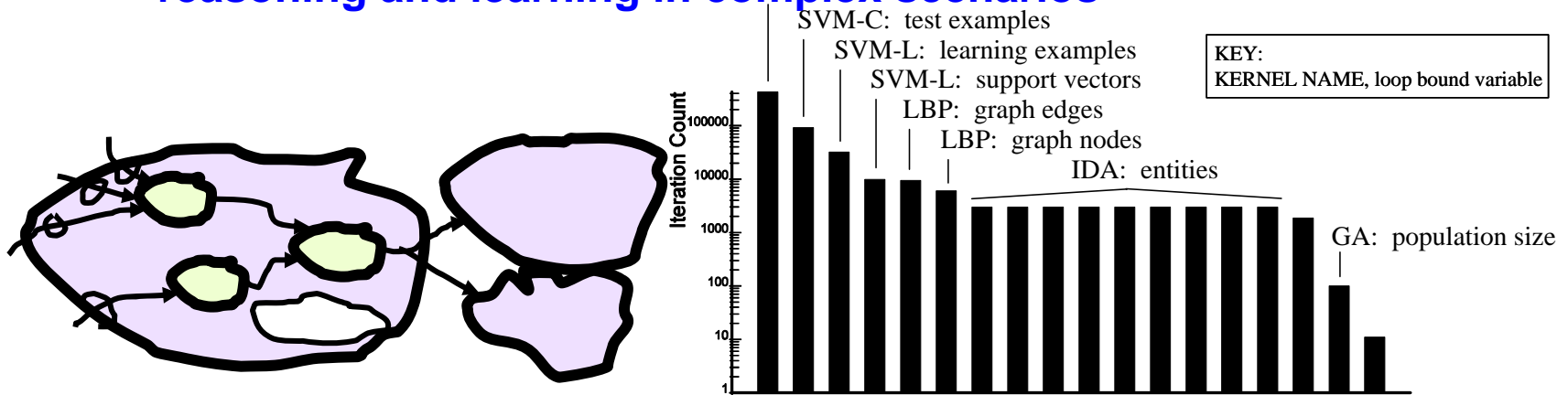
**Cell-based
introspection and
management**



Scalable Web of Cognitive Virtual Processing Elements



- Cognitive processing requires massive fine-grained parallelism with highly efficient processing elements
- Cognitive processing elements different from general-purpose computing, scientific computing, and signal processing elements
 - Processing granularity highly variable and dynamic
 - Cognitive systems and scenarios lead to dynamic code and data movement and load balancing
 - Density of parallelism must be much higher to do real-time reasoning and learning in complex scenarios



Parallelism With Varying Granularity and Computation Types

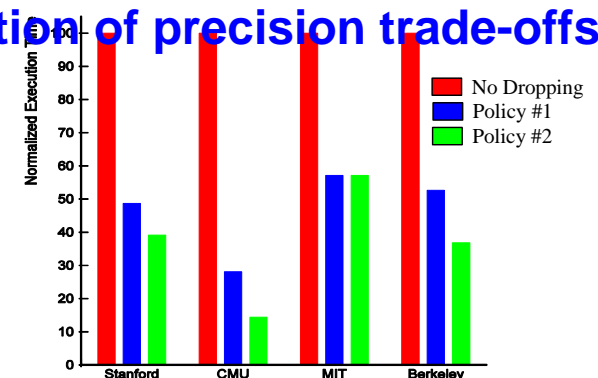


Multi-Level Soft Computing



- **Exploit the tolerance for imprecision, uncertainty, partial truth, and approximation to achieve tractability, robustness and low solution cost***
 - **Optimality or exactness infeasible in cognitive application domains**
 - **Input data has imprecision and inaccuracy**
 - **Robustness needed to handle transient and persistent faults**
- **Exploitation of soft computing for performance gains changes architecture at all levels**
 - **Processor: data types, functional units, circuit design**
 - **Memory: local and shared lossy memory protocols, latency reduction**
 - **Communication: lossy protocols, QoS tuning**
 - **System software: data types, communication of precision trade-offs to programmer, resource management**

Performance Improvements From Message Dropping



*<http://www.soft-computing.de/def.html>



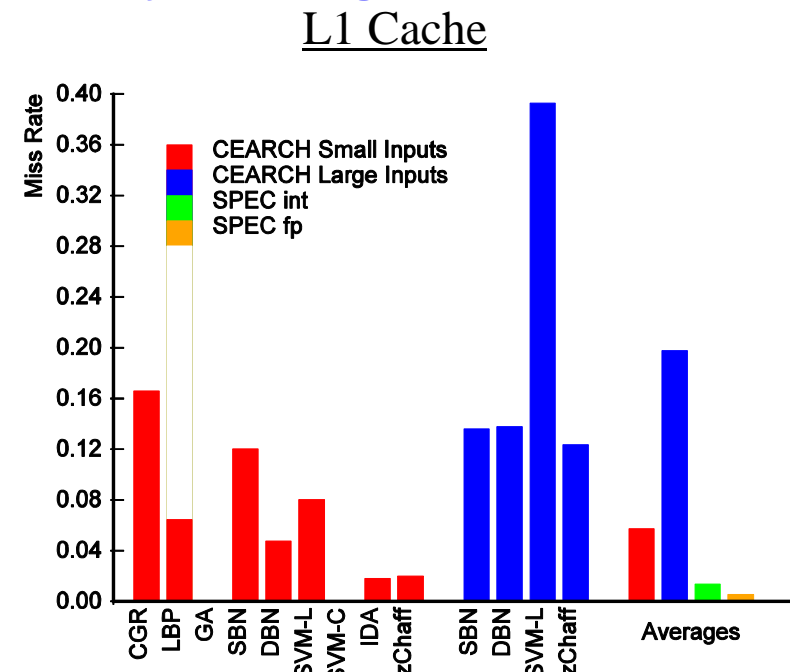
Adaptive Memory System



- Cognitive processing leads to poor memory system behavior in traditional memory systems
 - Some algorithms have irregular and hard-to-predict access patterns
 - Working sets can be very large because of complexity of scenarios
 - Dynamic resource allocation and fine-grained parallelism leads to more global memory accesses and locality challenges

- Memory system requirements

- Flexible allocation among cognitive processing elements
- Fine-grained protection
- Flexible commit policies
- Inexpensive roll-back for fault tolerance and race conditions between parallel compute elements



Miss Rates for Cognitive Algorithms Using Traditional Cache



CEARCH Architecture Layers



Programming Model

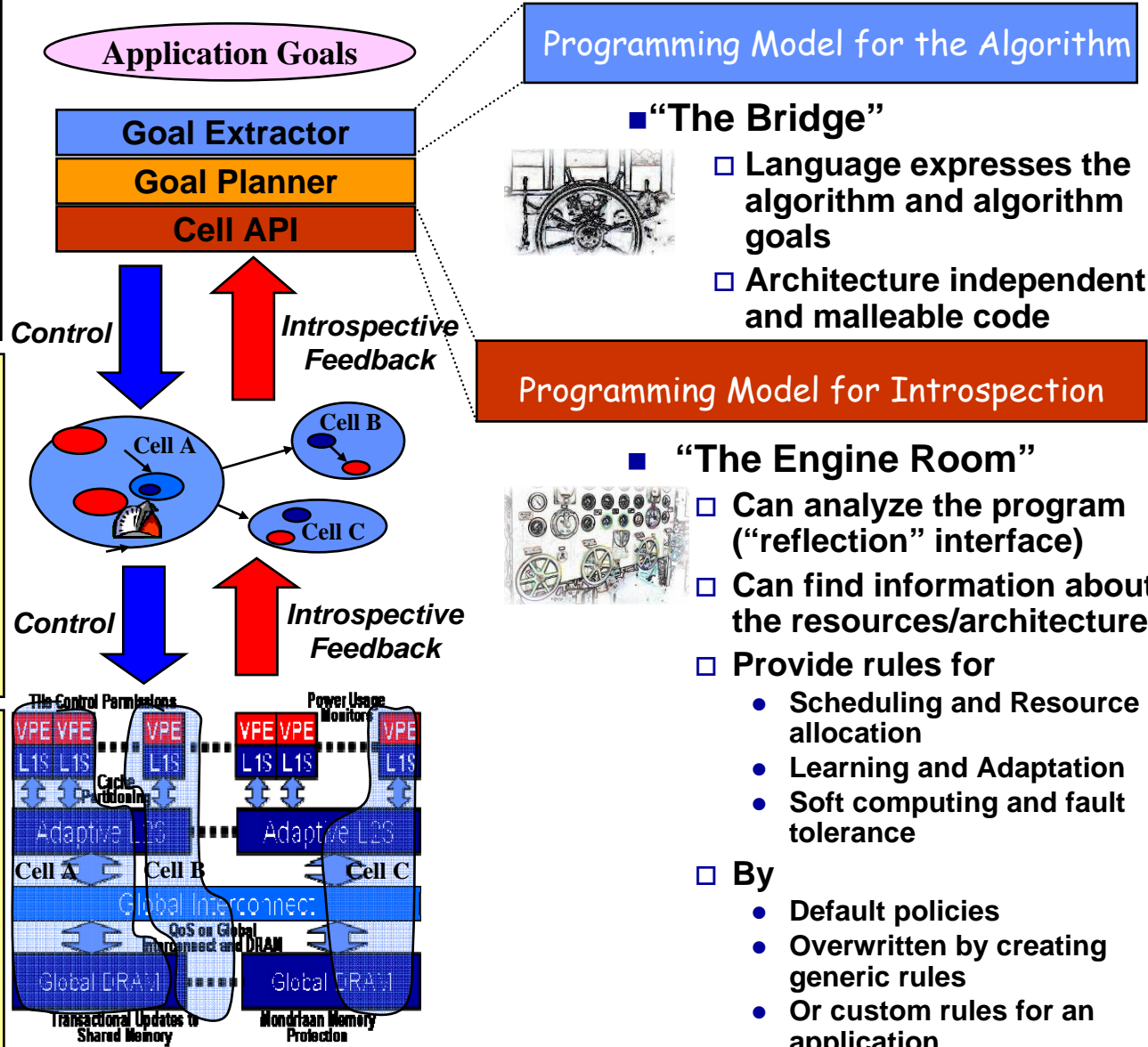
- Abstraction barriers provide scalable low-level performance with high-level specifications
- Goal-based performance and resource allocation allows computation to be in part selected by system
- Soft computing semantics

Runtime System

- Learning and reasoning-based goal-oriented instrumentation and compilation
- Adaptive and introspective hierarchical resource allocation for processing, memory, and communication

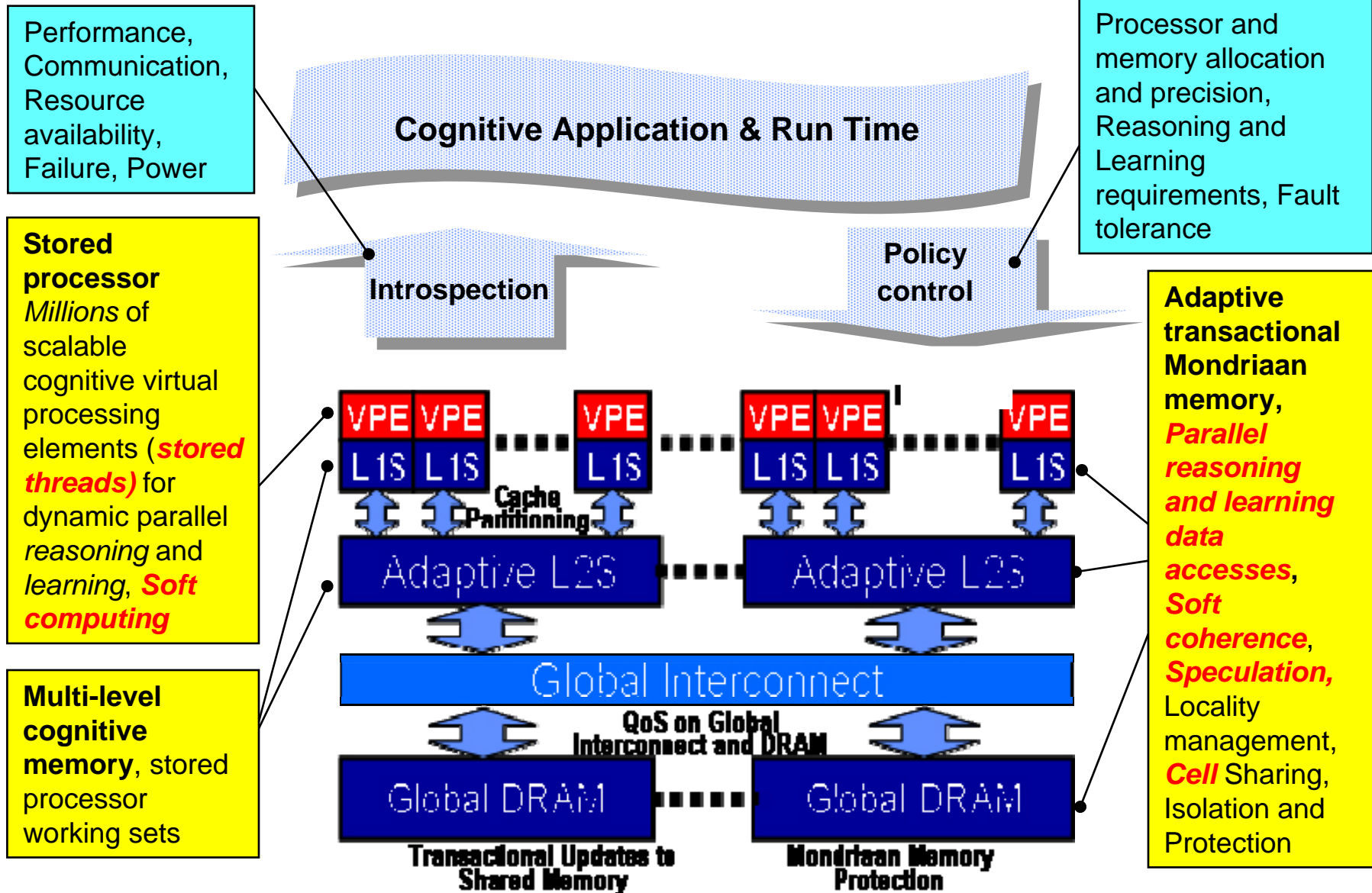
Hardware Architecture

- Millions of introspective virtual processing elements running on thousands of hardware engines
- Adaptive memory for efficient data access and sharing
- Soft computing support





CEARCH Hardware Architecture





Outline



- **Project Goals**
- **Architecture Characteristics**
- **Application Examples**
- **Summary**



Spotting Behaviors OODA Loop

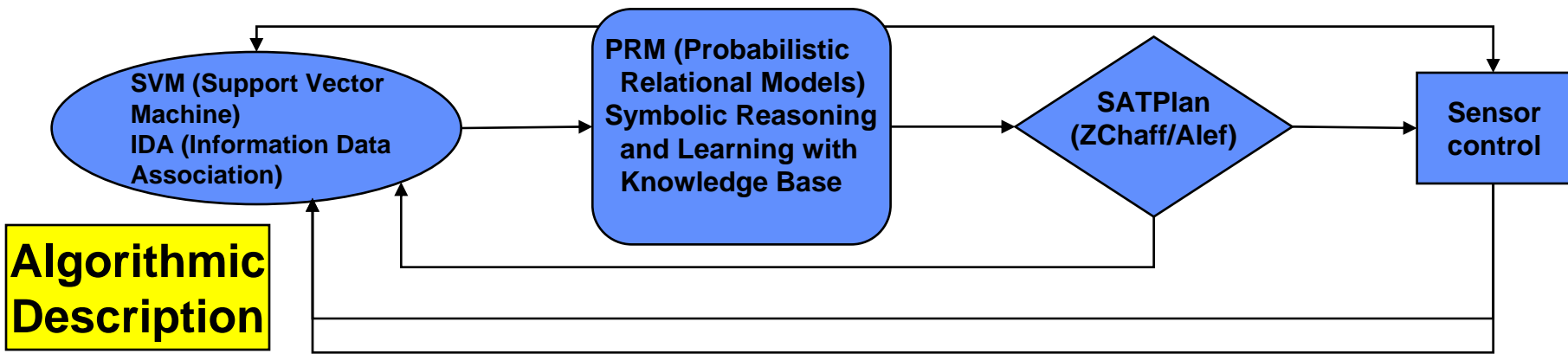
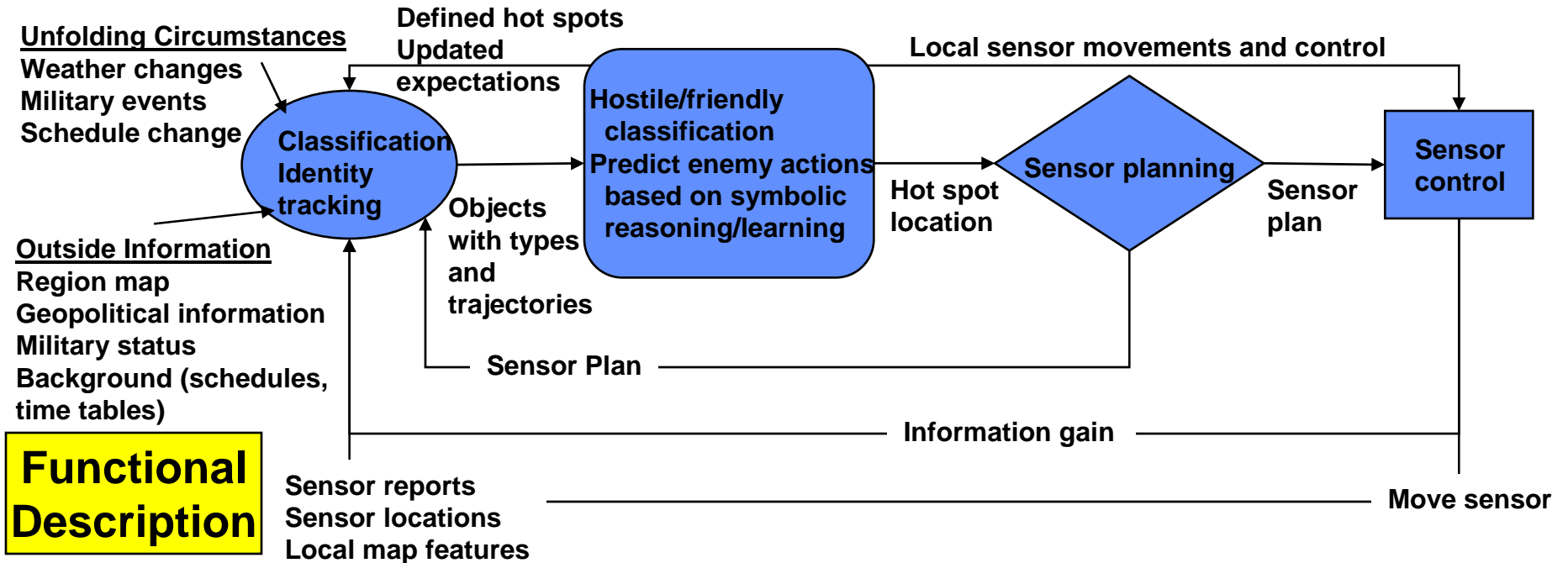


Observe

Orient

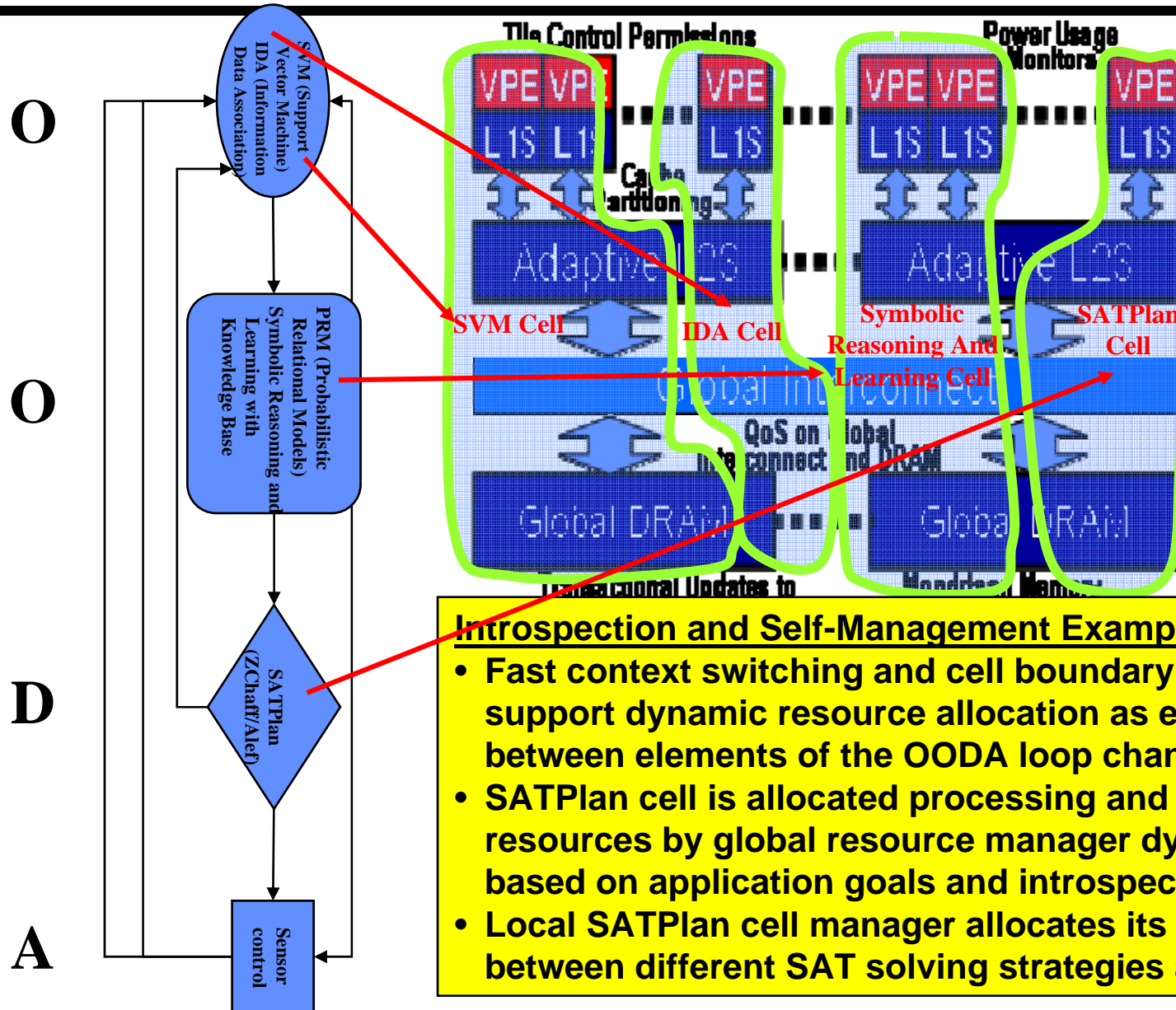
Decide

Act





Introspection and Self-Management

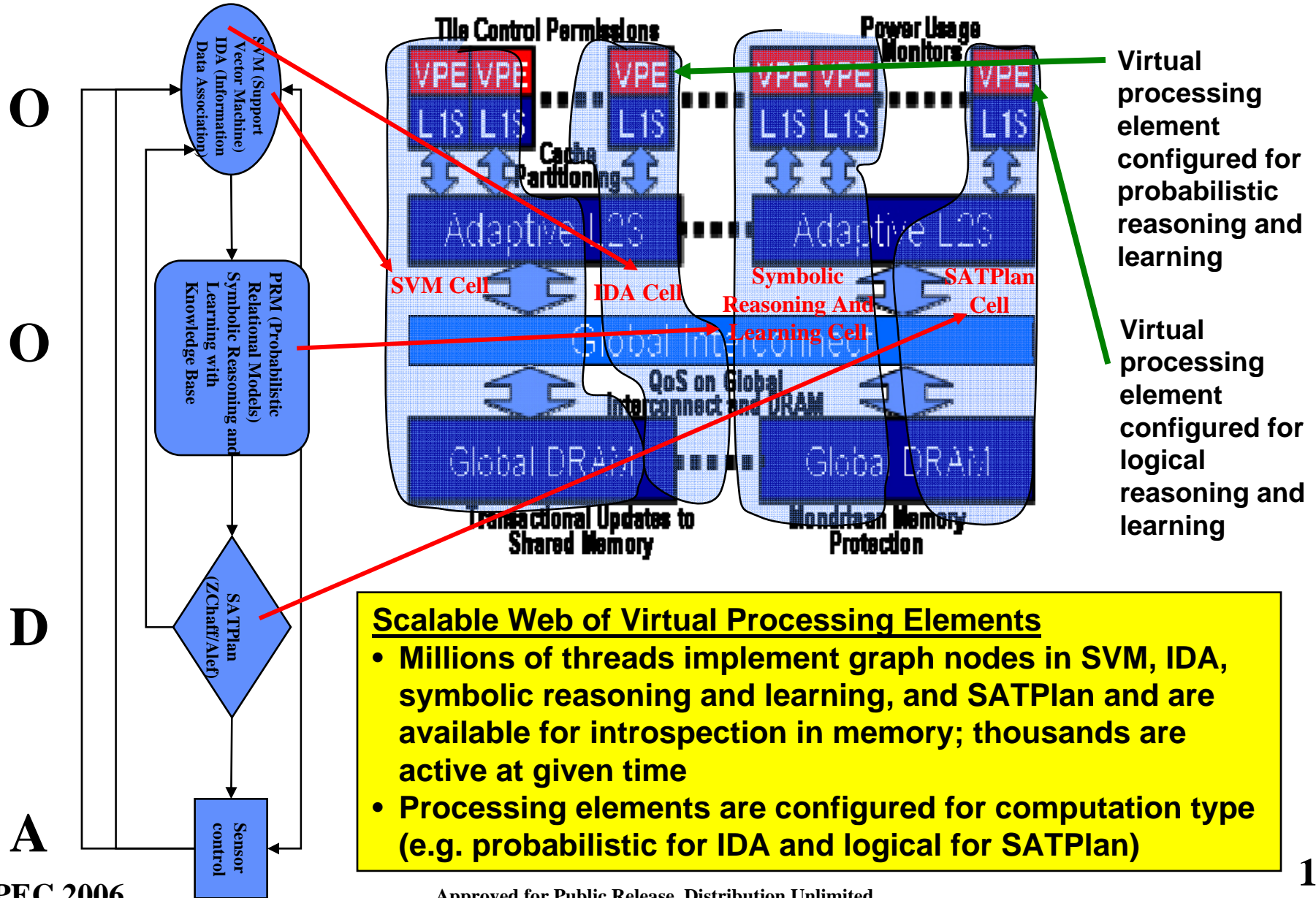


Introspection and Self-Management Examples

- Fast context switching and cell boundary changes support dynamic resource allocation as emphasis between elements of the OODA loop changes
- SATPlan cell is allocated processing and memory resources by global resource manager dynamically based on application goals and introspective monitors
- Local SATPlan cell manager allocates its resources between different SAT solving strategies and sub-goals

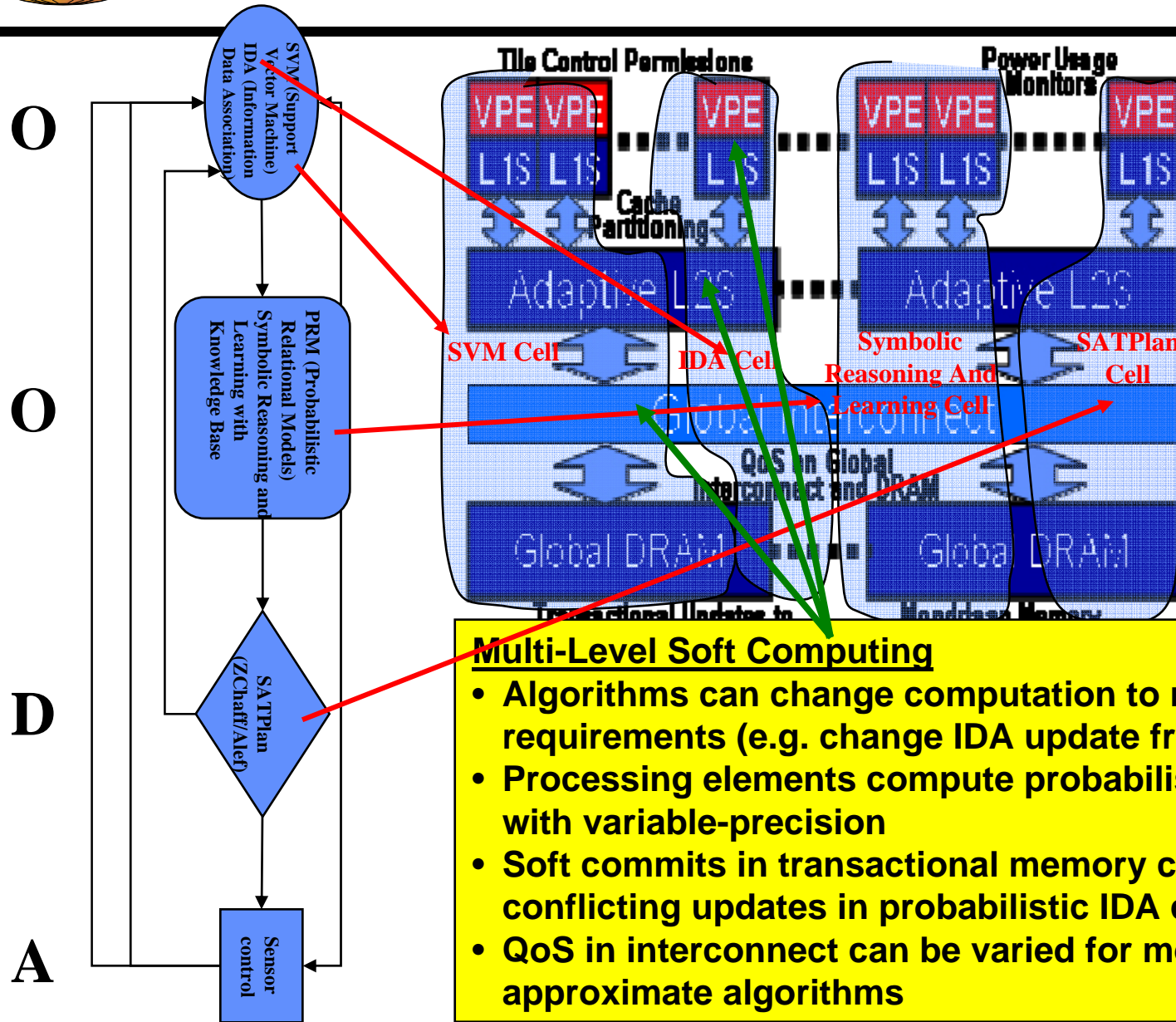


Scalable Web of Virtual Processing Elements



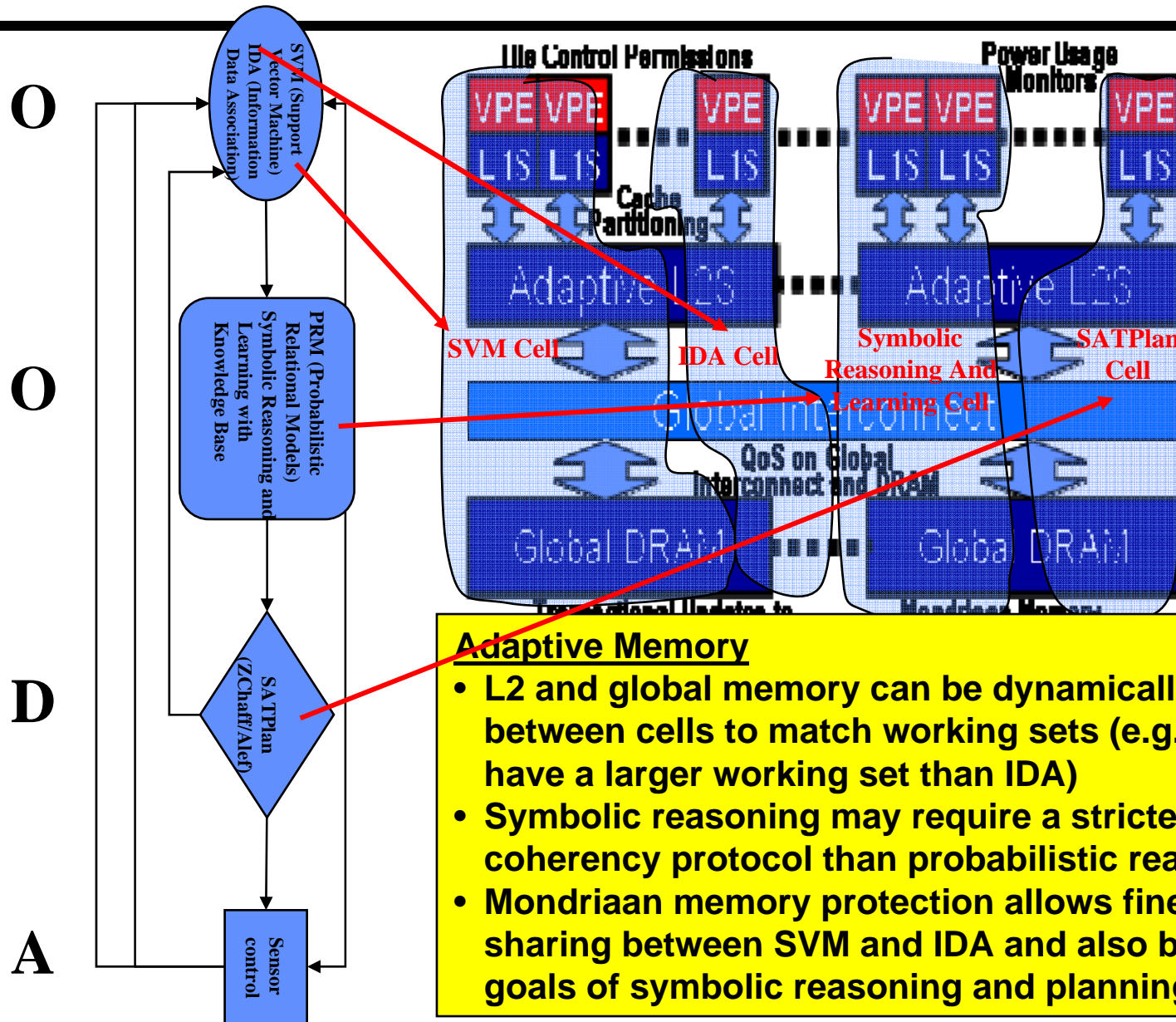


Multi-Level Soft Computing



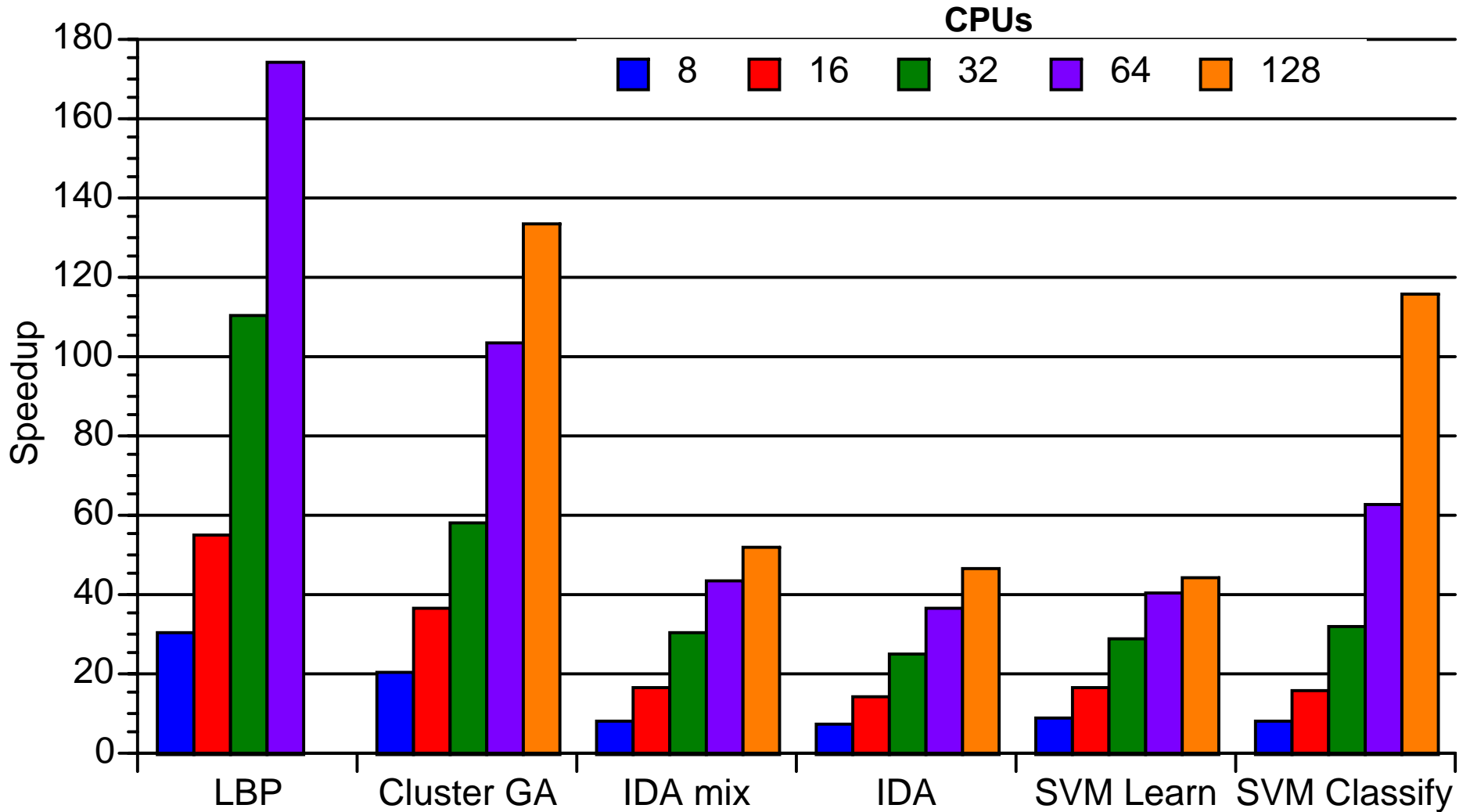


Adaptive Memory



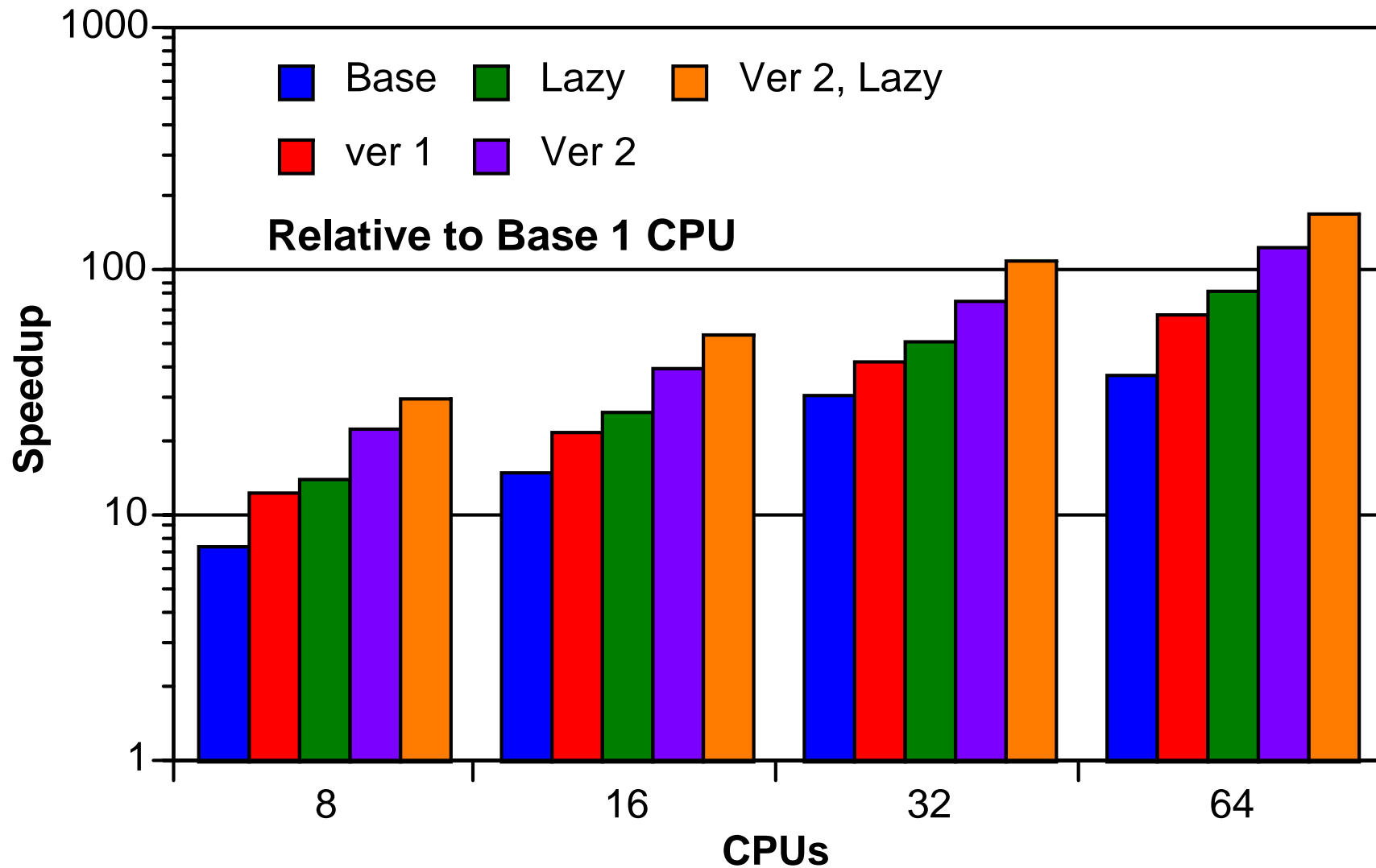


CEARCH Application Speedups





LBP Performance Improvement





Outline



- **Project Goals**
- **Architecture Characteristics**
- **Application Examples**
- **Summary**



CEARCH Summary



- **CEARCH is a dynamic self-managing architecture for cognitive processing uniquely suited to complex environments**
 - Driven by cognitive system and algorithm characteristics
 - Dynamically organize resources to optimize performance, power and reliability
 - Adaptation and introspection in both hardware and software
- **CEARCH has unique features to efficiently support cognitive applications and that provide capability not possible with today's COTS architectures**
 - Stored processor
 - Adaptive, transactional memory
 - Soft computation
 - Introspection and run-time policy control support
- **Preliminary architecture evaluation indicates**
 - High performance potential
 - Well suited to cognitive applications and soft computing