# X-Stack, Traleika Glacier (XSTG)

Three-year project results and products, including Intel, Rice, UCSD, PNNL, UIUC, Reservoir Labs, & Eqware

**Shekhar Borkar, Josh Fryman**

**&**

**The TG Team**

# Outline

Overview of extreme-scale programs

Architecture & software stack

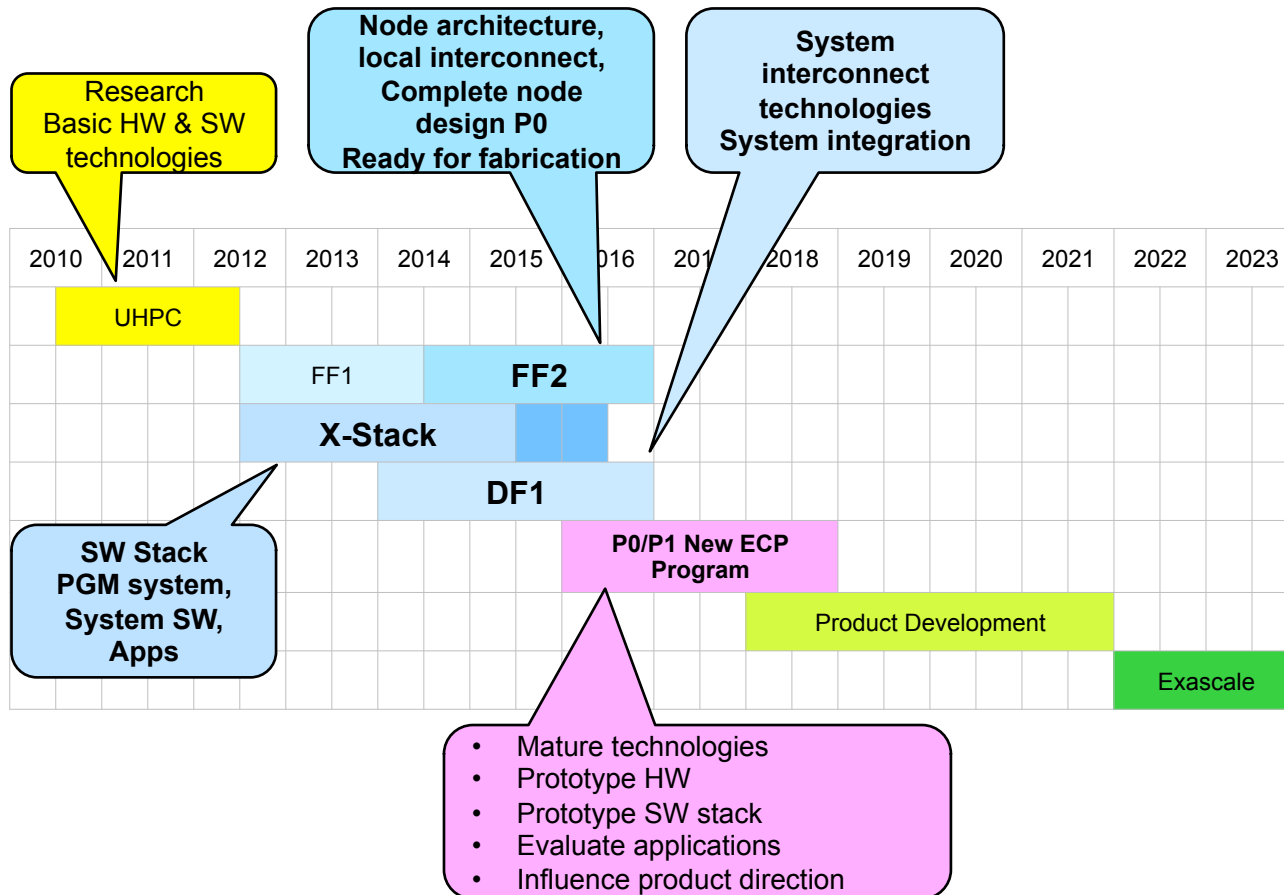Results of Tools and Applications
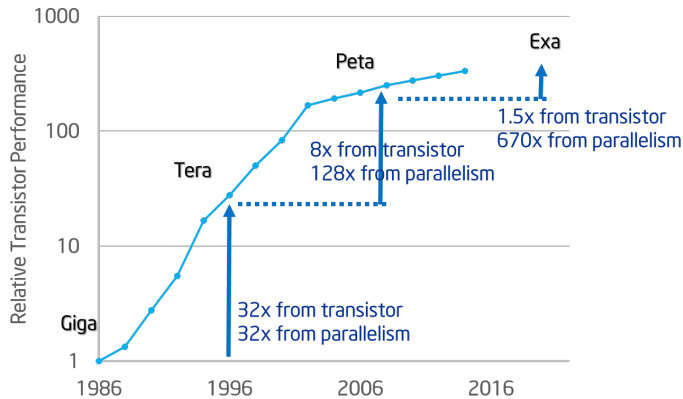
Community Responses

Report Card

Future Work

# Top Extreme-scale Challenges

1. **System Power & Energy** **(Exascale in 20 MW)**

2. **New, efficient, memory subsystem** **(BW, latency, capacity, cost)**

3. **Extreme parallelism O(Billion)** **(productivity, legacy)**

4. **New execution model** **(introspection, event-driven)**

5. **Resiliency to provide system reliability** **(>> 6 days)**
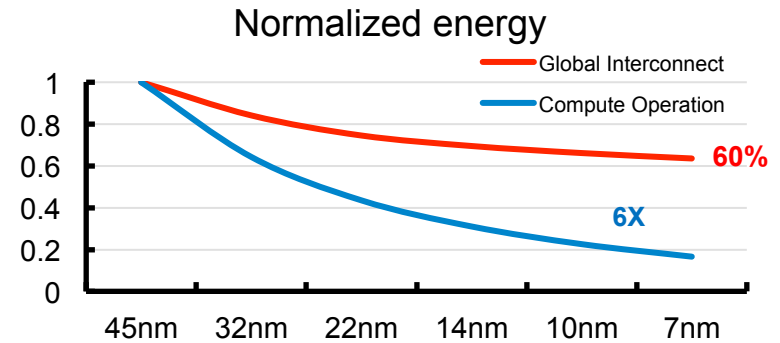
6. **$ Cost and affordability by enhancing system-efficiency**
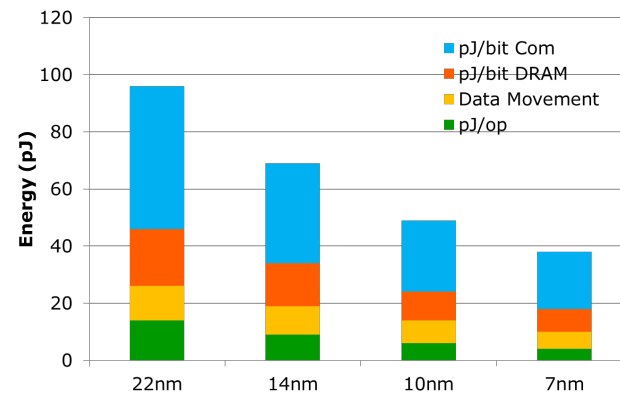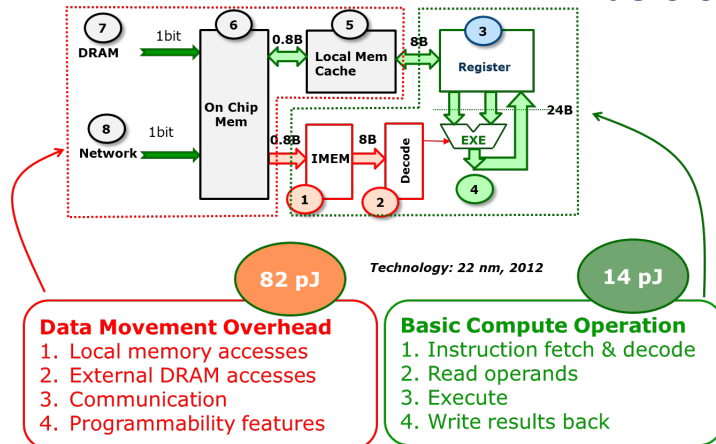
# Extreme-Scale Programs for <u>Ubiquity</u>



**Research**
Basic HW & SW technologies

**Node architecture, local interconnect, Complete node design P0 Ready for fabrication**

**System interconnect technologies System integration**

| 2010 | 2011 | 2012 | 2013 | 2014 | 2015 | 2016 | 2017 | 2018 | 2019 | 2020 | 2021 | 2022 | 2023 |
|------|------|------|------|------|------|------|------|------|------|------|------|------|------|

UHPC

FF1   FF2

X-Stack

DF1

SW Stack PGM system, System SW, Apps

P0/P1 New ECP Program

Product Development

Exascale

- Mature technologies
- Prototype HW
- Prototype SW stack
- Evaluate applications
- Influence product direction

4

# Compute vs Data Movement



**Performance from parallelism**

1000 — Exa — 1.5x from transistor / 670x from parallelism

Peta — 8x from transistor / 128x from parallelism

Tera — 32x from transistor / 32x from parallelism

Relative Transistor Performance: 1986, 1996, 2006, 2016

## Normalized energy



Global Interconnect — 60%
Compute Operation — 6X

45nm 32nm 22nm 14nm 10nm 7nm

## Basic compute loop



Technology: 22 nm, 2012

**82 pJ**

**Data Movement Overhead**
1. Local memory accesses
2. External DRAM accesses
3. Communication
4. Programmability features

**14 pJ**

**Basic Compute Operation**
1. Instruction fetch & decode
2. Read operands
3. Execute
4. Write results back



- pJ/bit Com
- pJ/bit DRAM
- Data Movement
- pJ/op
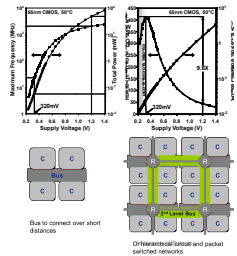
Energy (pJ): 22nm, 14nm, 10nm, 7nm
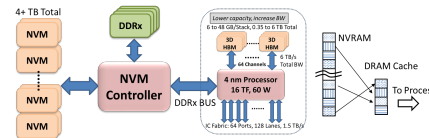
5

# Memory Hierarchy



6

# Addressing Challenges
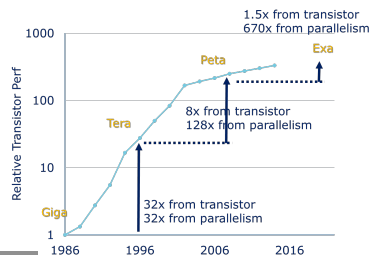
## 1. Power & Energy

Vdd scaling

NTV operation

Fine grain power management

Hierarchical, heterogeneous IC fabric

## 2. Memory subsystem

Disproportionately large local memories

Incorporate NVM for cost

Hierarchy, global addressability

## 3. Programmability

O(B) threads

Data locality

Separation of concerns

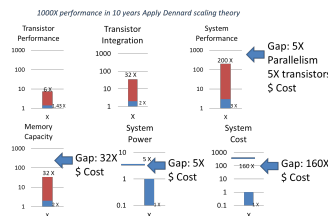Legacy compatibility

Productivity

## 4. Execution model

Asynchronous

Event-driven tasks (tiny threads)

Dynamic scheduling of threads

Introspective resource management

(Self-aware through observations)

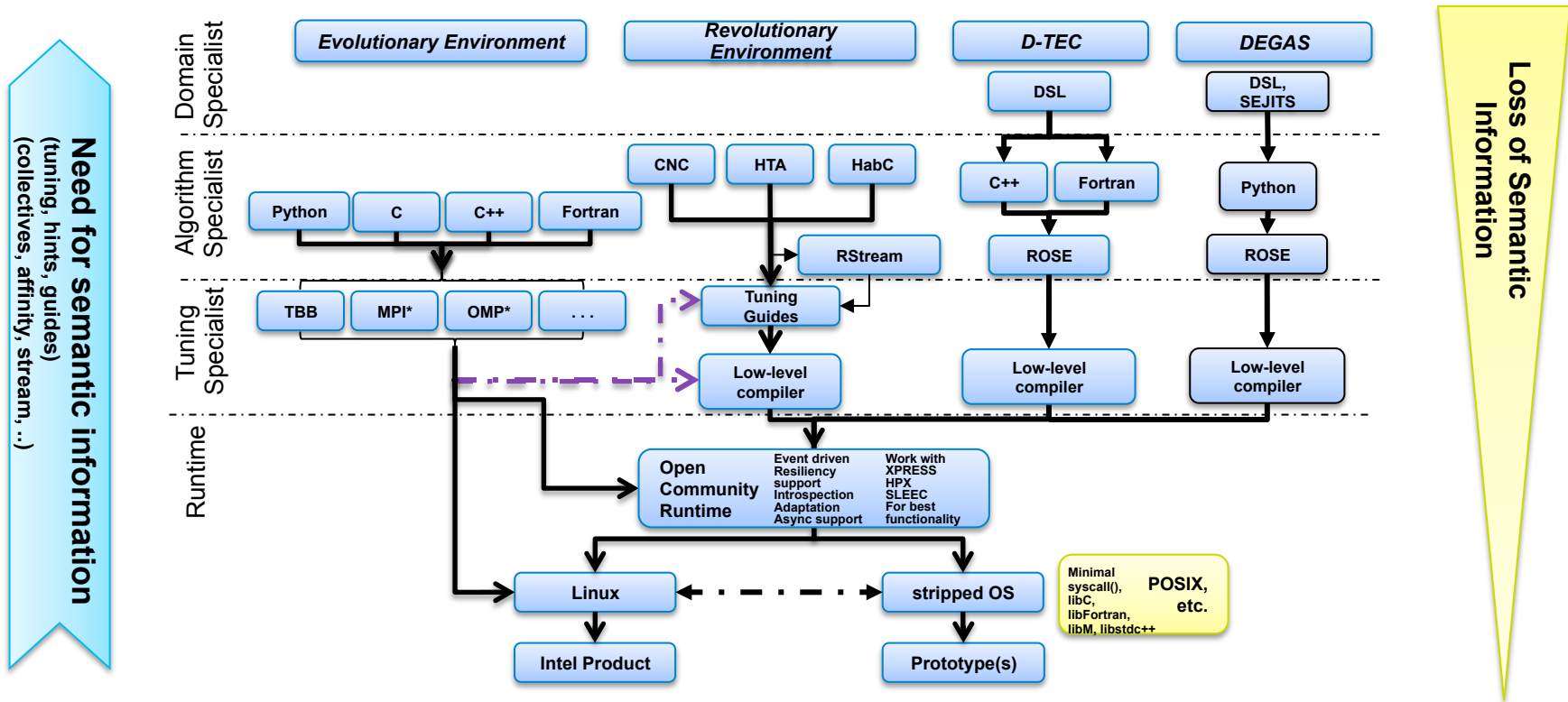Dynamic system optimization

## 5. Resiliency

First, understand faults

Detection in HW with sensors

Reactive and proactive measures

## 6. Cost & affordability

Improve Perf/Transistor

New, efficient architecture

Data locality

Insert NVM in the hierarchy

# Software Ecosystem: Support today, find tomorrow

# Major Points

- All work done under XSTG is available publicly, open-source BSD licensed

  - http://xstack.exascale-tech.com

  - Exceptions:
    - PNNL has not yet released their OCR implementation at this time
    - R-Stream is closed-source with government use rights

- XSTG project has created significant reference implementations

  - Formal specification for task-based execution model

  - Reference implementation that even in immature form exceeds expectations

  - Productivity tools, compilers, profilers, visualizers, re-factored applications, libraries of kernels and examples, simulators, and tutorials

  - **Application workshops every 6 months with the community**

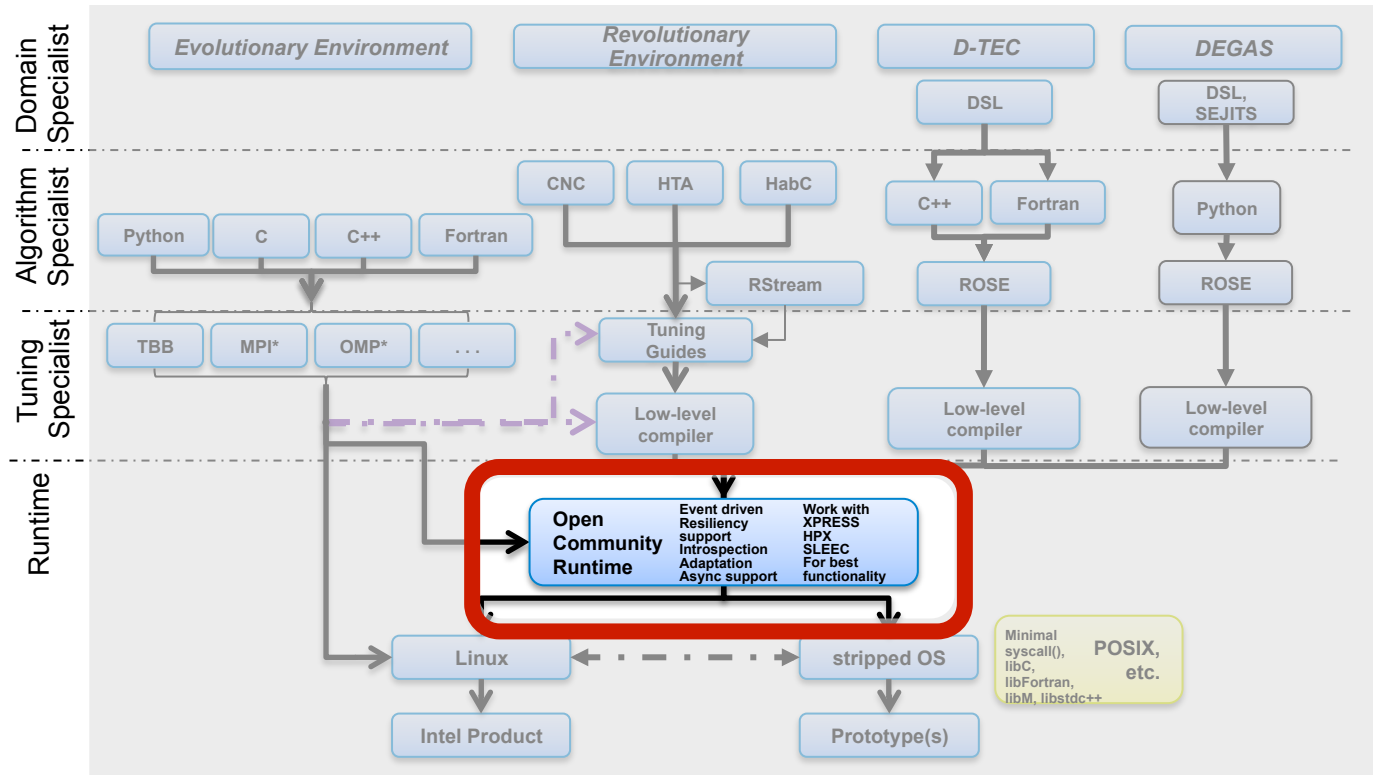    *XSTG has transitioned from Research to Technology Maturation*

# Embargoed Data

- The following sections contain results that have not yet been published

- Please embargo sharing and distribution of this content until publication efforts have concluded by each team

- All results are identified by team – contact each for further details

# Task-Execution Runtime (Intel, Rice)
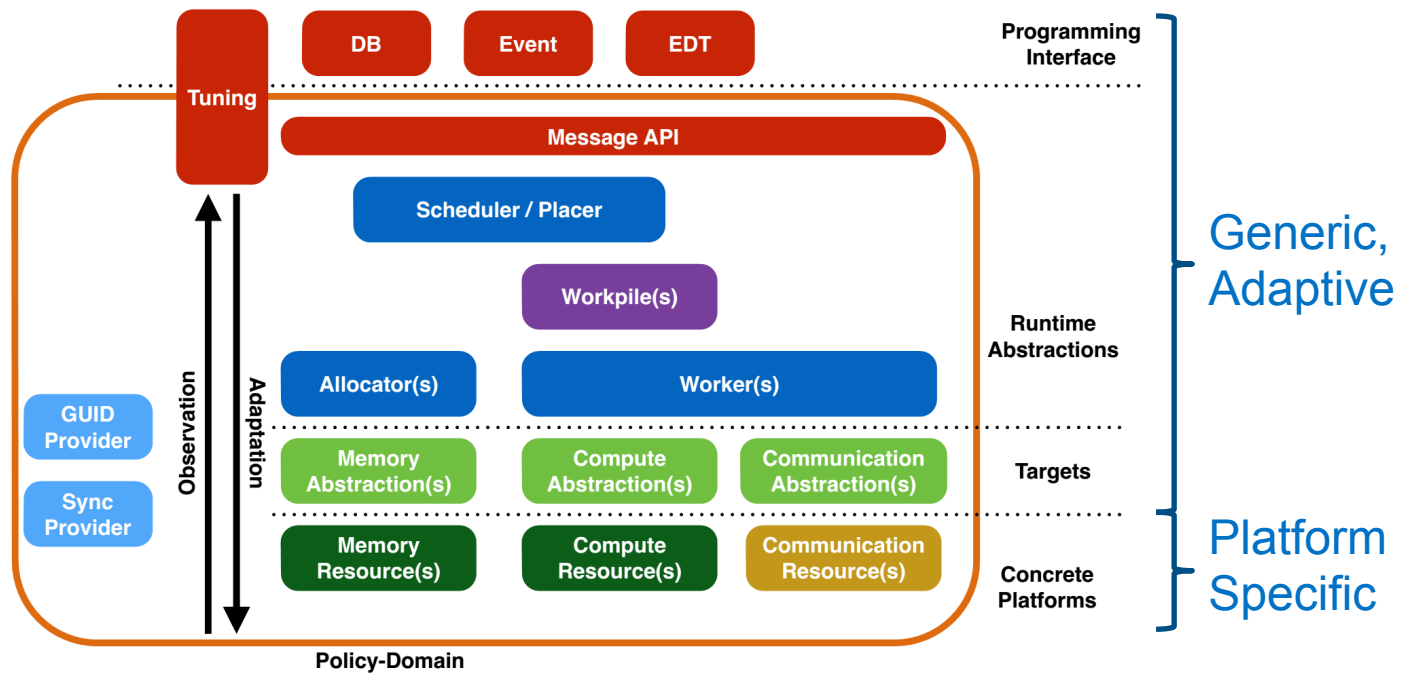
# Open Community Runtime (OCR)

## Multi-party collaboration

- Intel, Rice, UIUC, UCSD, PNNL, Reservoir Labs

- Provide effective abstraction for diverse hardware (hetero-ISA ready)

- Typify future task-based execution models

- Handle large-scale parallelism efficiently and dynamically

- Provide user-perspective application-transparent resiliency

- Maintain a separation of concerns (application/scheduling/resources)

- Open source (encourage collaboration) http://xstack.exascale-tech.com

- *OCR is X-Stack Traleika Glacier project's implementation for this revolutionary run-time prototype*
- *FFWD-2 is extending OCR with legacy support, re-factored applications, and re-factoring guides, templates, and tools*
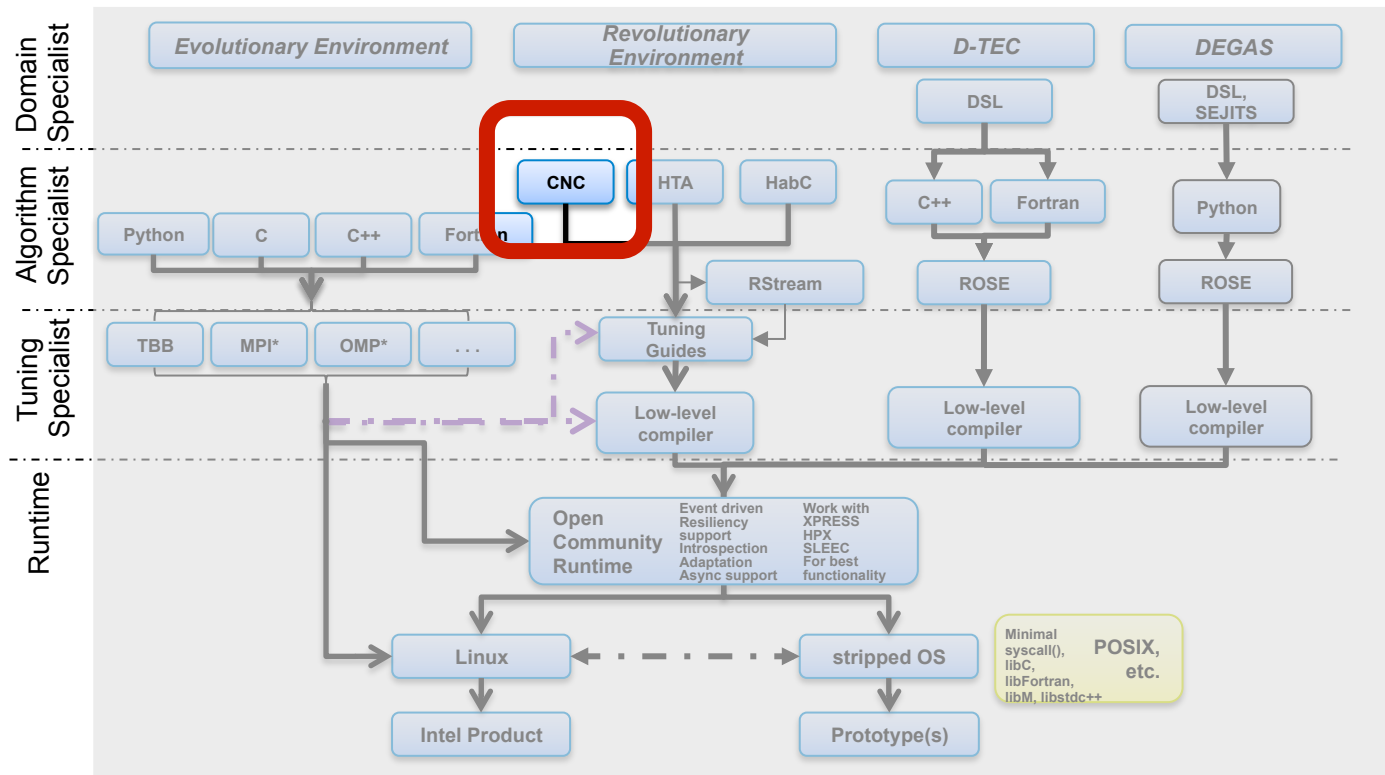
# Runtime Design Principles of OCR



Modular

Extensible

Adaptable

Tunable

Generic, Adaptive

Platform Specific

# Productivity: Concurrent Collections (Rice, Intel)

# Productivity: Hierarchically Tiled Arrays (UIUC)

# Productivity: R-Stream (Reservoir Labs)

# R-Stream: Polymorphic Optimizing Translator

```
for(k=1;k<NPAD-1;k++){
    for(j=1;j<NPAD-1;j++){
        for(i=1;i<NPAD-1;i++){
            x_np1_1[k][j][i] = x_np1_0[k][j][i]
                            c1*(x_np1_0[k][j][i]
                            c2*Dinv_ijk()* (rhs[
}}}}}
```

Input C code

- Automatic parallelization and locality optimization
- Automatic explicit communications generation and explicit datablock management
- Automatic formation of OCR EDT programs (and other runtimes)
- Dynamic creation and management of dependences



Non-trivial loop boundaries and array access functions in the core of computation block of each EDT

Optimized R-Stream-- OCR code embeds non-trivial dependence management for dynamic creation of EDTs and datablocks

R-Stream outputted OCR code snippet

# Task-Execution Runtime (PNNL)

# Software Ecosystem: Support today, find tomorrow

# OCR and SoC Methodology: Co-Design

# Execution Engine (XE) architecture – tunable to target

# DOE Community Responses and Thoughts (1/2)
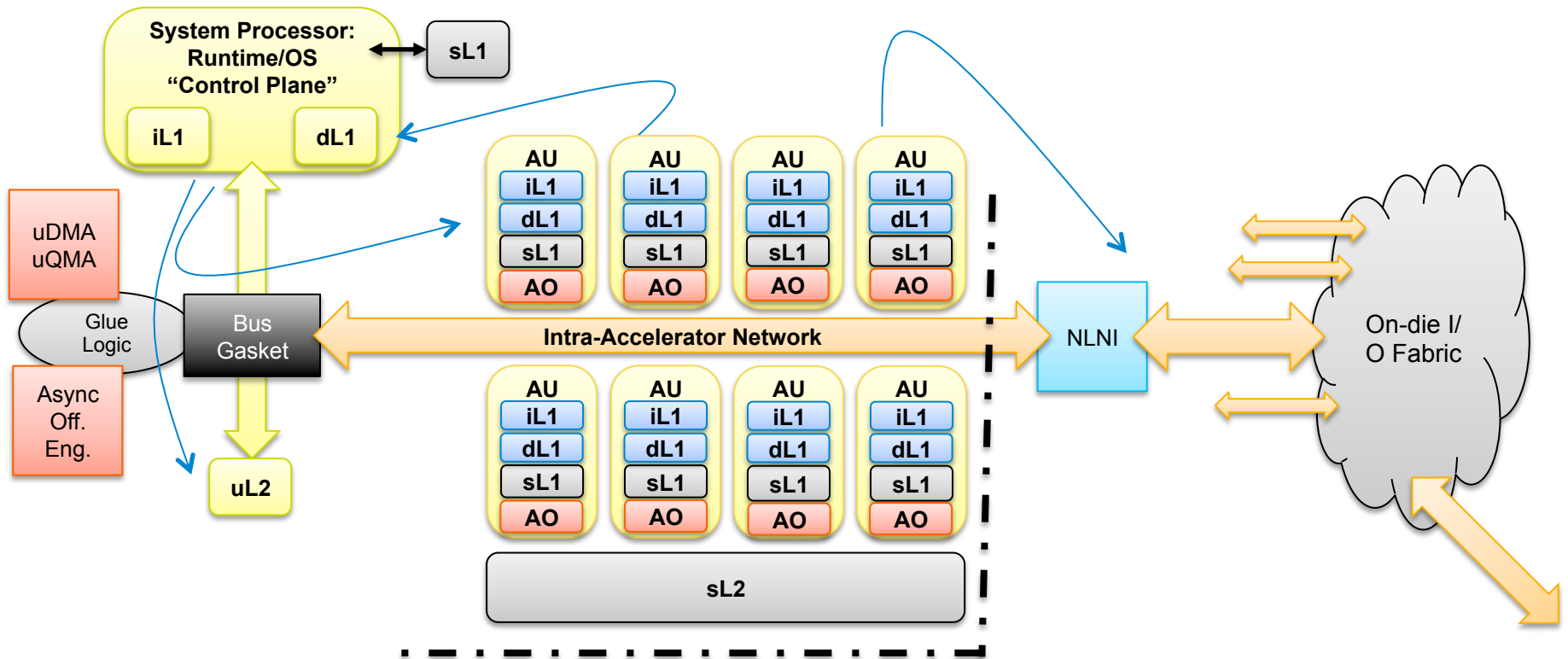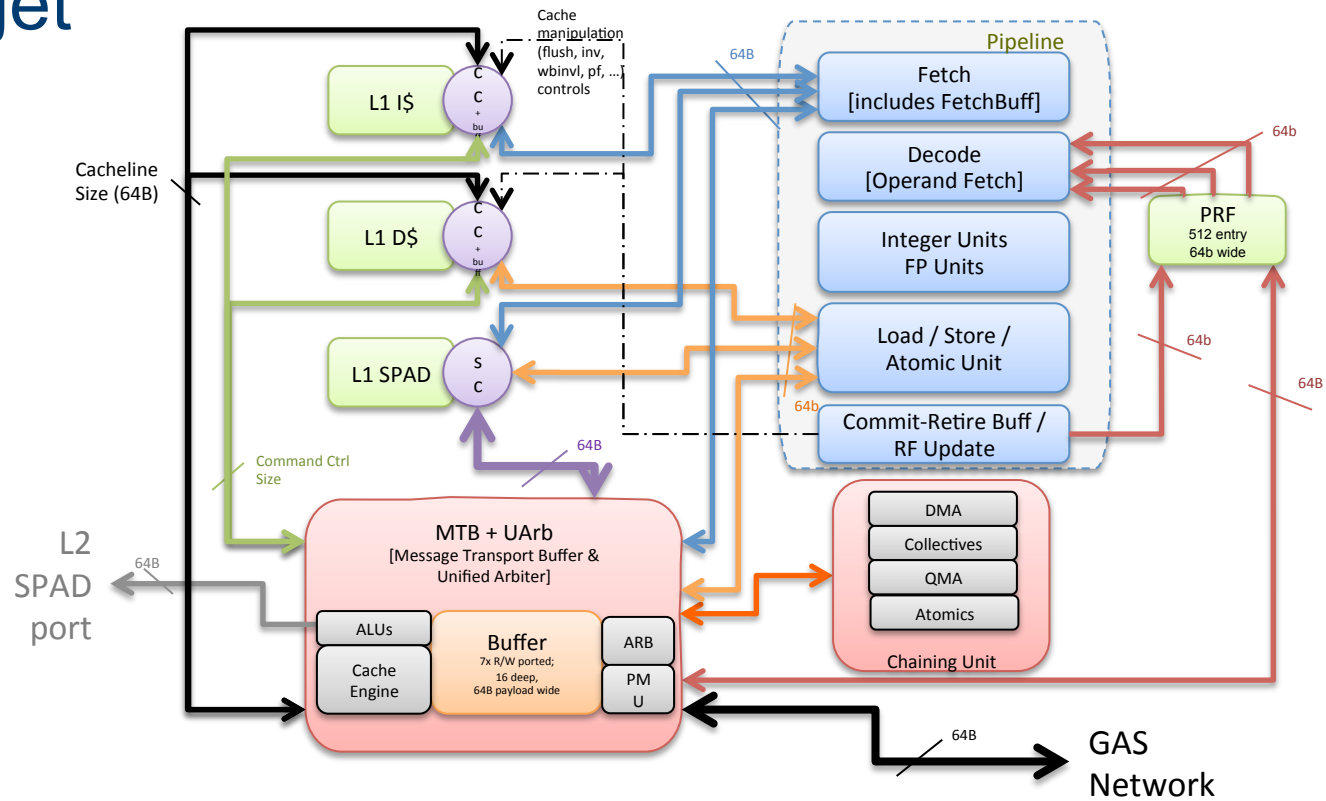
| Member | DOE Lab | Comments |
|--------|---------|----------|
| Dave Richards | LLNL | I'm **pleased by the progress of the OCR project.** MPI-Lite, a full specification, and the expanded feature set including affinity hints are important additions that will expand the set of problems OCR can handle. My biggest concern with OCR is that we have yet to see a believable story regarding how high-level programming models can interface with OCR and capture the full benefits of asynchronous tasks and relocatable data blocks. |
| Pat McCormick | LANL | The Legion programming model from LANL will be moving on top of OCR to assess performance and viability. If successful, **Legion will use OCR as one of the primary back-end runtimes** to provide task-based execution model support. |
| Robert Clay | SNL-L | I'm hugely supportive of the approach — specification-based, open design and development. This is what the community needs to build interchangeable 'components' in the task-parallel RTS space. **We (Sandia) plan to remain actively engaged with the OCR team(s), and consider OCR an extremely interesting candidate** for a core part of our future task-parallel RTS in context of a DHARMA build out. While OCR doesn't seem ready for prime time use today, I believe with continued community engagement and support it will develop into a highly effective capability which will significantly enhance its adoption w/in the broader community. |

# DOE Community Responses and Thoughts (2/2)

| Member | DOE Lab | Comments |
|--------|---------|----------|
| Ian Karlin | LLNL | The OCR research has two advantages over the other mainly university led asynchronous multi-task (AMT) runtimes. **OCR development is tied to hardware exploration, which allows for hardware features that are beneficial** to AMT runtimes to be co-designed with the programming model. In addition, it is working to become an open community product as opposed to a tool owned by one research group. That said OCR still suffers from the main drawback of most AMT models, as it has yet to be proven that adopting an AMT model is significantly beneficial for enough codes over a well designed non-bulk synchronous MPI application. |
| Jim Belak | LLNL | I am very pleased by the overall progress and direction of the XSTG project. **I am also excited about the direction and vision that this project is exploring**. With several application codes now expressed in the model, analysis has revealed the need to introduce concepts, such as hints, to enable performant applications using the model. I look forward to seeing large, at-scale applications running to determine where and when this approach is most appropriate for our problems |

# Progress Report Card

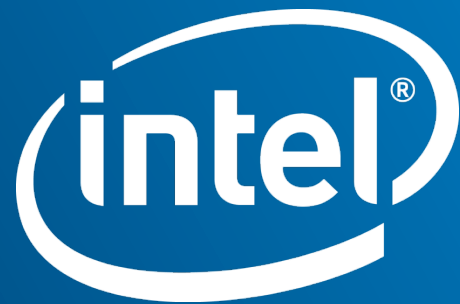| Technology | Rating | Comments |
|---|---|---|
| CnC | Good | CnC specification and reference are open-sourced by Intel; Rice is now driving changes and enhancements |
| HTA | Stopped | HTA works, but is relatively immature in the ability to express constructs and realize performance on OCR. Further work at UIUC is needed on the basic design. |
| Habanero | Good | Habanero C support for OCR works well. The lack of a centralized standard or reference implementation for Habanero hampers ability to support OCR with all Habanero features and versions. |
| R-Stream | Good | R-Stream is very effective at automatically generating high-performance OCR code from high-level C loop nests, that exploit the full feature set of OCR and allow users to get the full benefit of the EDT model vs. classical execution models. |
| OCR (Intel, Rice) | Good | Functionally complete with a precise open specification and reference implementation. Advanced features on resiliency, hints, and related technologies are in progress. |
| OCR (PNNL) | Good | A version derived from the open reference implementation with a focus on performance first. Source code is not available. |
| FSim | Good | Very fast, highly parallel, large-scale architecture simulator works and matches the RTL design model at a functional level. Supports bare-metal as well as OCR applications. |
| Tools (debug, profile) | Limited | First-pass tools are working, but are not going to scale as-is to the expected size of future machines without significant additional work. Collaboration with DOE groups working on similar tools is just starting, and will remain in progress. Starting this was delayed until the last year of the XSTG project, due to the evolving OCR design and interfaces. |

# Future Work: 1 of 2

| Technology | Comments |
|---|---|
| Maturation - Xeon | Mature code base on scheduler and allocator to bring up performance and strengthen design corners |
| Maturation – Xeon Phi | Study KNL-based architecture optimization opportunities for performance and introspection |
| Maturation – FSim | With RTL design closing on specific performance characteristics and power values, adjust FSim to reflect and optimize the tools targeting FSim and OCR on FSim accordingly |

# Future Work: 2 of 2

| Technology | Comments |
|---|---|
| Productivity - Legion | Support DOE task-based framework Legion on top of OCR to provide a comparison baseline |
| Productivity - RAJA | Exploit loop lambda model of RAJA to emit OCR tasks for parallelization evaluation; couple to Rstream and similar tools for "auto-taskification" skeleton for existing codes |
| Enhancements - Tools | Continue to enhance and extend early debug and profiling tools with DOE teams to facilitate using task-based execution models at scale |

© 2015 Intel Federal LLC