



center for
cognitive computing
systems research

Codesigning Cognitive Computing Systems and Applications

Wen-mei Hwu, Co-Director with Jinjun Xiong (IBM)

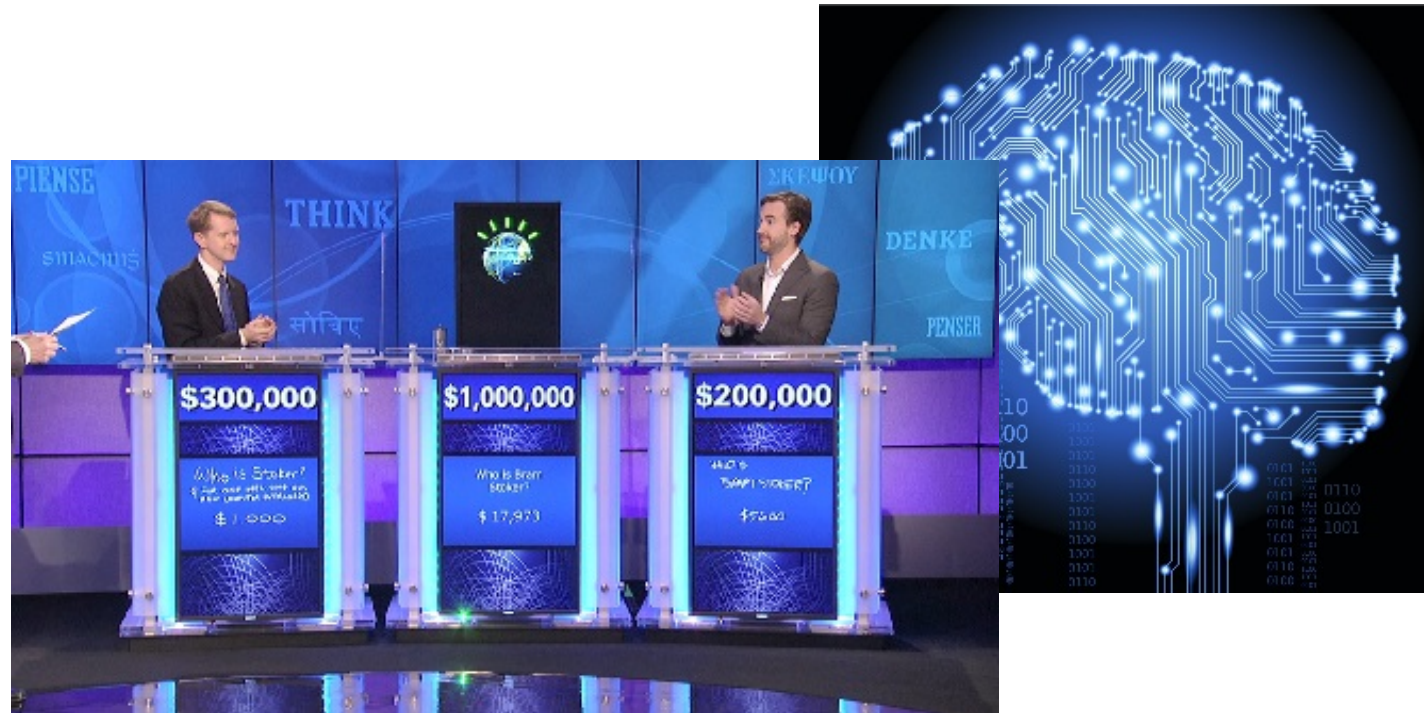
on behalf of the entire C³SR team

University of Illinois at Urbana-Champaign

November 2, 2017

Cognitive Computing – the C3SR View

- A cognitive computing application fuses vast, **unstructured data** and vast **human knowledge base** to **extend human capabilities** by solving problems, making actionable recommendations, and producing customized learning experiences



C3SR Vision

- The rise of cognitive computing has created new opportunities to rethink all the three layers of computing systems—applications, software, and hardware.
- Dramatic enhancement in the **efficacy**, **efficiency** and **variety** of cognitive computing applications can be achieved through dramatic enhancement in the **programmability**, **throughput**, **latency**, **capacity**, and **affinity** of computing systems.

C³SR faculties & students (Est. 9/2016)



Suma Bhat



Minh Do



Deming Chen



Julia Hockenmaier



Wen-mei Hwu



Nam Sung Kim



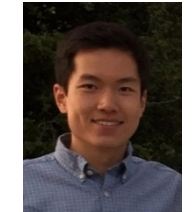
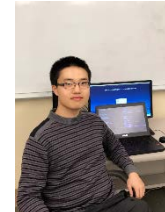
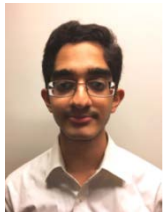
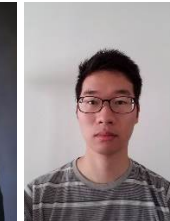
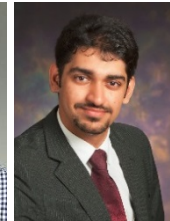
Dan Roth



Rakesh Nagi



Lav Varshney



...

The Three Pillars of C3SR:

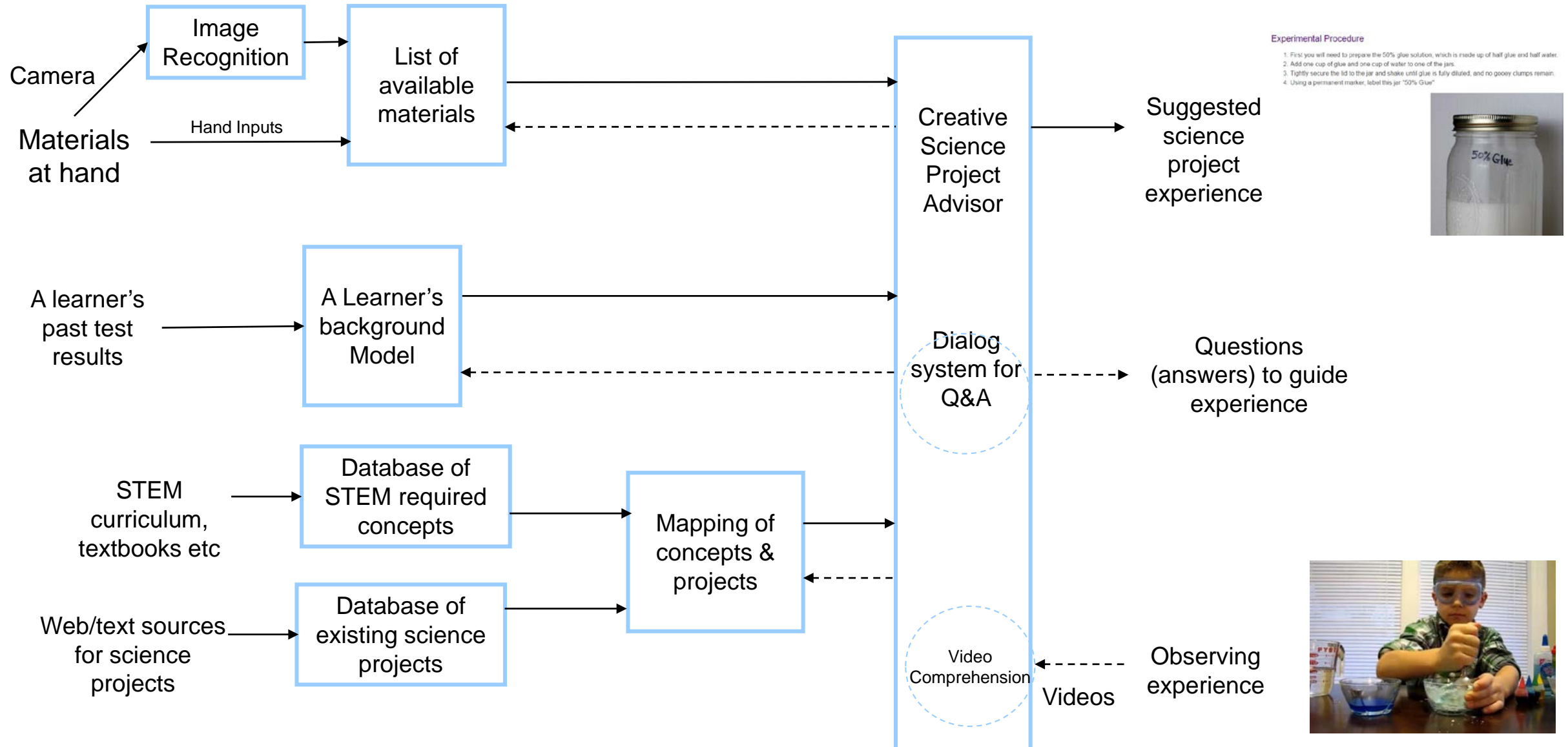
- Creative experiential learning advisor (CELA) as a grand challenge use case for cognitive capabilities
- Cognitive application builder (CAB) to make the underlying heterogeneous infrastructure easy to consume for cognitive application developers
- Cognitive systems innovations (Erudite) for workload acceleration, including Near Memory Acceleration (NMA)

A New Modality of Application Development

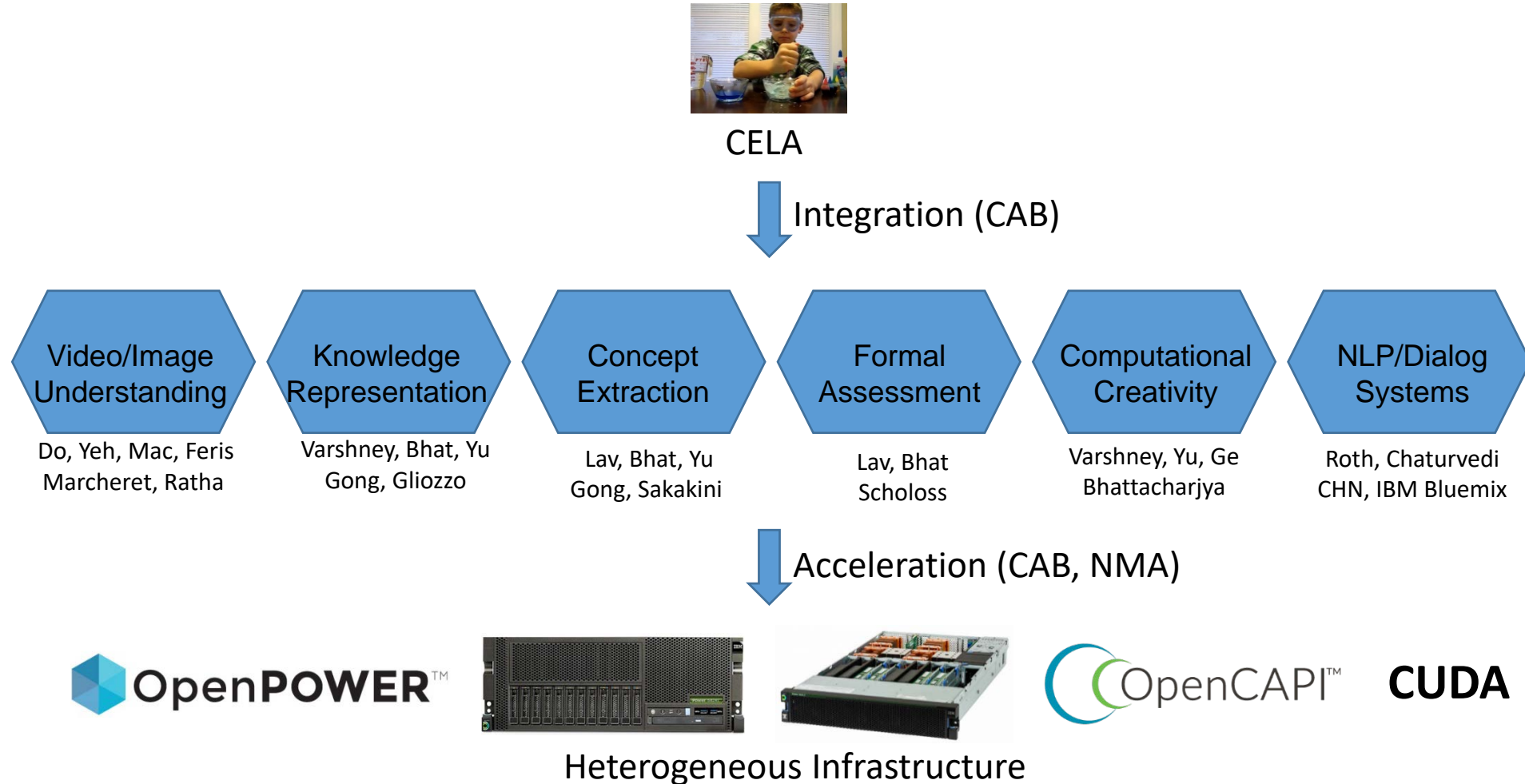
- Cognitive applications demand functionalities that we have failed to program
 - Computer vision, natural language dialogs, stock trading, fraud detection, ...
- Use labeled data – data that come with the input values and their desired output values – to learn what the logic should be
 - Capture each labeled data item by adjusting the program logic
 - Learn by example!
- This introduces a new modality of application development
 - Training, Testing, Integration, Profiling, Debugging, etc.

Application Driver

CELA: personalized education via multi-modality data comprehension and computational creativity



Decomposition of CELA's Research Challenges



- Requires a tool to integrate core services that are optimized for the underlying heterogeneous infrastructure

Cognitive Application Builder

A system-level challenge

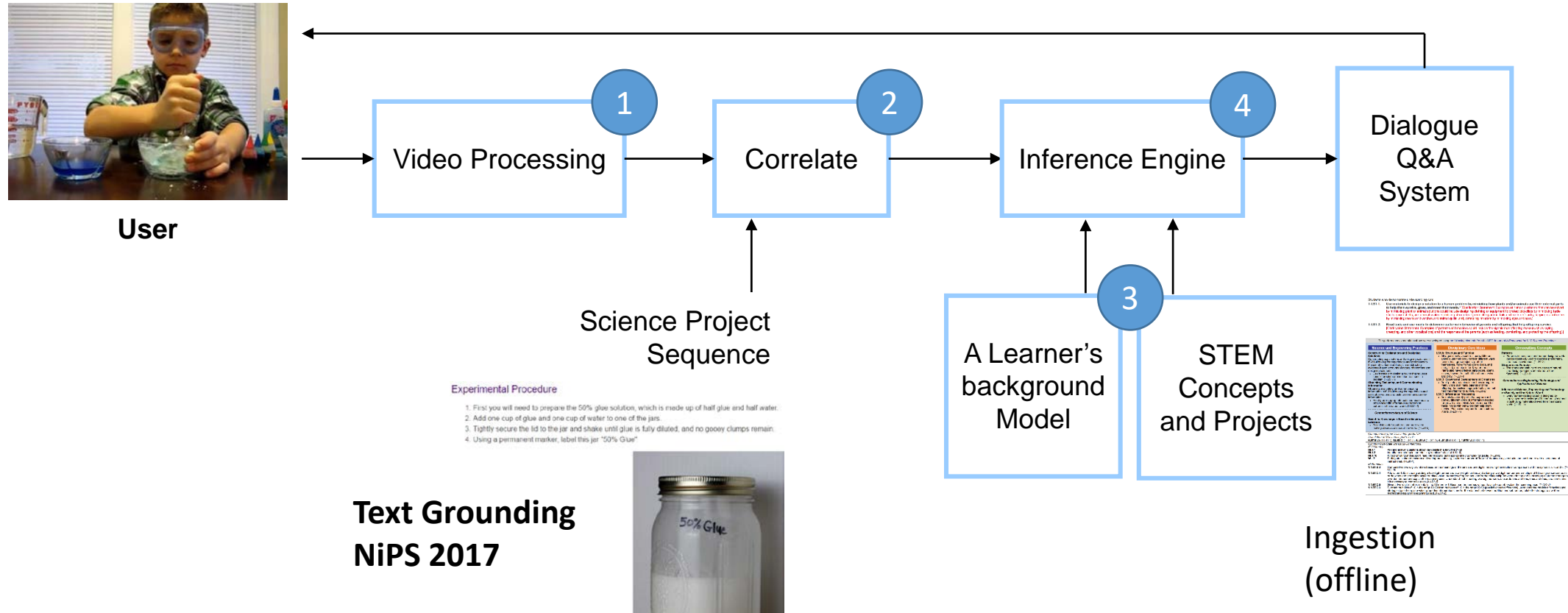
Workflow description
Innovative AI techniques,
Data, Models, Frameworks



High-performance, scalable,
robust applications

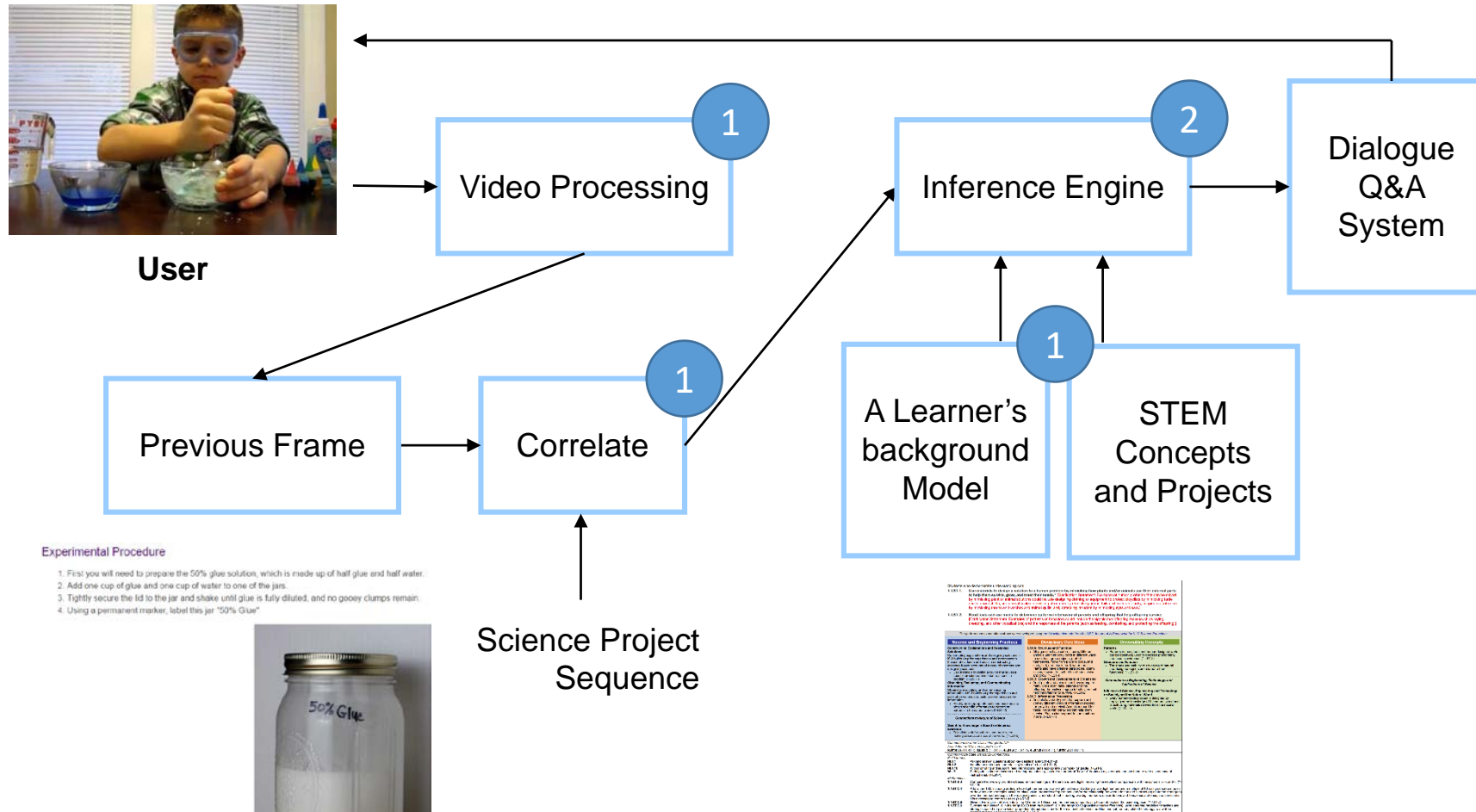
- CAB: A language, compiler, and runtime for easy development of cognitive applications
 - Software synthesis to exploit accelerators and efficient communication
 - Introspection for debugging and performance evaluation
 - Workflow profiling, optimization and orchestration for system-level performance
 - Decentralized application architecture for scalability, composability, testing, and development

CELA as a Driving Use Case for CAB



CAB simplifies component connection, workflow description, model training/selection, and iterative development

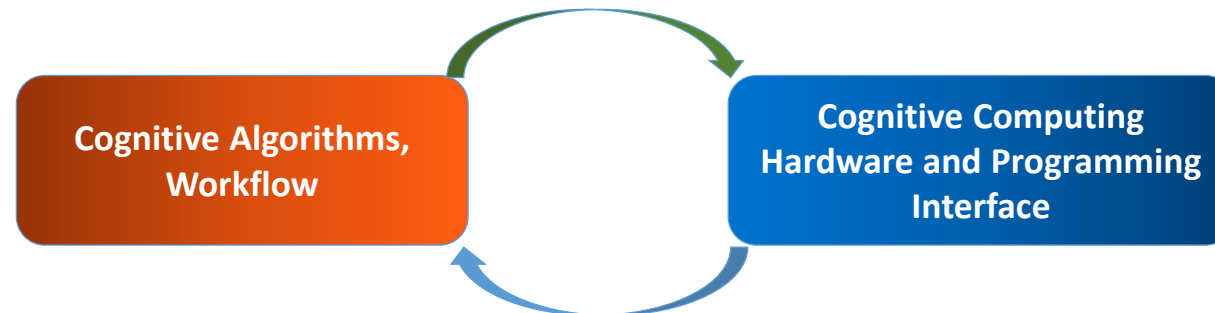
CELA as a Driving Use Case for CAB



- CAB will automatically transform workflows for high-performance execution

C3SR Approach to Cognitive Computing System Design

- To develop scalable cognitive applications by co-designing
 - advanced methods and algorithms for cognitive computation, and
 - optimized heterogeneous computing systems for these workloads.
- Generations of complete prototype systems
 - **Initial** – existing methods, algorithms and workflows running on existing hardware
 - **Refined** – innovative methods, algorithms and workflows enabled by the next generation memory/storage technology and accelerators
 - **Novel** – ambitious methods, algorithms and workflows empowered by new memory and near memory/near IO acceleration technologies.



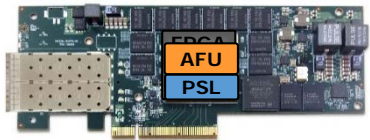
Initial Experimental Heterogeneous Infrastructure



2 x P8 Minsky with NVLink GPUs



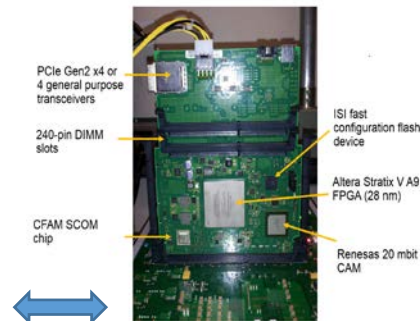
DGX-1



FPGA CAPI
over PCIe



4 x P8 Tuleta (S824L)



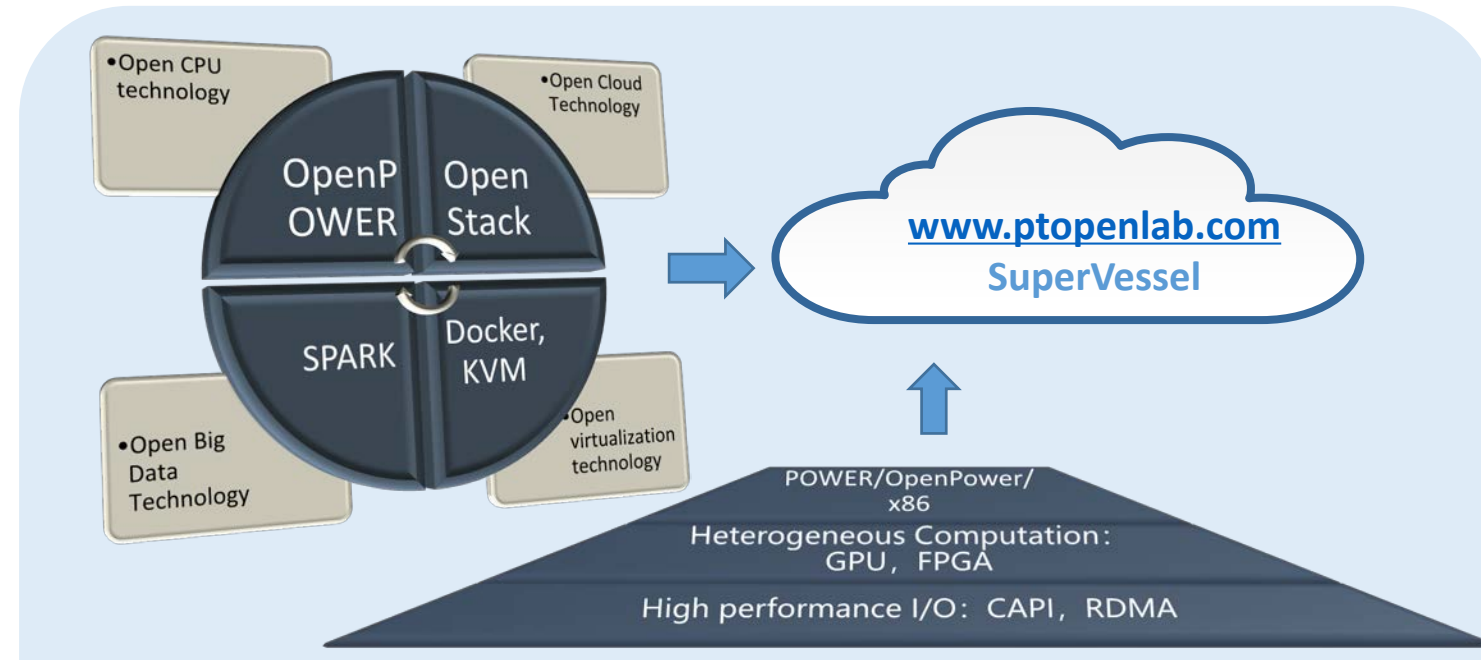
ConTutto over DMI



IBM Bluemix™



Watson developer cloud



Great support from Brad McCredie

- A dedicated program manager and team (Ben Kreuz, JT Kellingon, Adam McPadden, Dean Sciacca, Jonanthan Dement)

Selected center progress highlights

- Curated datasets**
- The CarML system for model development and deployment**
- Workload acceleration**
- The Erudite NMA system**

Curated Datasets

- Extracted STEM concept dependency from next generation science standard that includes

- Performance Expectations
- Science and Engineering Practices
- Disciplinary Core Ideas
- Crosscutting Concepts
- Connections

Students who demonstrate understanding can:

1-LS1-1. Use materials to design a solution to a human problem by mimicking how plants and/or animals use their external parts to help them survive, grow, and meet their needs.* [Clarification Statement: Examples of human problems that can be solved by mimicking plant or animal solutions could include designing clothing or equipment to protect bicyclists by mimicking turtle shells, acorn shells, and animal scales; establishing structures by mimicking animal tails and roots on plants; keeping out intruders by mimicking thorns on branches and animal quills; and, detecting intruders by mimicking eyes and ears.]

1-LS1-2. Read texts and use media to determine patterns in behavior of parents and offspring that help offspring survive. [Clarification Statement: Examples of patterns of behaviors could include the signals that offspring make (such as crying, sleeping, and other vocalizations) and the responses of the parents (such as feeding, comforting, and protecting the offspring).]

The performance expectations above were developed using the following elements from the NRC document *A Framework for K-12 Science Education*:

| Science and Engineering Practices | Disciplinary Core Ideas | Crosscutting Concepts |
|--|---|--|
| Constructing Explanations and Designing Solutions Constructing explanations and designing solutions in K-2 builds on prior experiences and progresses to the use of evidence and ideas in constructing evidence-based accounts of natural phenomena and designing solutions. • Use materials to design a device that solves a specific problem or a solution to a specific problem. (1-LS1-1) Obtaining, Evaluating, and Communicating Information Obtaining, evaluating, and communicating information in K-2 builds on prior experiences and uses observations and texts to communicate new information. • Read grade-appropriate texts and use media to obtain scientific information to determine patterns in the natural world. (1-LS1-2) ----- Connections to Nature of Science Scientific Knowledge is Based on Empirical Evidence • Scientists look for patterns and order when making observations about the world. (1-LS1-2) ----- Connections to other DCIs in first grade: N/A. Articulation of DCIs across grade levels: 1-LS1-1 (1-LS1-1); 2-LS2-2 (1-LS1-1); 4-LS1-1 (1-LS1-1); 4-LS1-2 (1-LS1-1); 4-ETS1-1 (1-LS1-1) Common Core State Standards Connections: ELA/Literacy RI.1.1 Ask and answer questions about key details in a text. (1-LS1-2) RI.1.2 Identify the main topic and relevant key details of a text. (1-LS1-2) RI.1.3 With prompting and support, read informational texts appropriately complex for grade. (1-LS1-2) RI.1.4 Participate in shared research and writing projects (e.g., explore a number of "How-to" books on a given topic and use them to write a sequence of instructions). (1-LS1-1) Mathematics 1.NBT.B.3 Compare two two-digit numbers based on the meanings of the tens and one digits, recording the results of comparisons with the symbols >, =, and <. (1-LS1-2) 1.NBT.C.4 Add within 100, including adding a two-digit number and a one-digit number and adding a two-digit number and a multiple of 10, using concrete models or drawings and strategies based on place value, properties of operations, and/or the relationship between addition and subtraction; relate the strategy to a written method and explain the reasoning used. (1-LS1-2) 1.NBT.C.5 Subtract within 100, including subtracting a two-digit number and a one-digit number and subtracting a two-digit number and a multiple of 10, using concrete models or drawings and strategies based on place value, properties of operations, and/or the relationship between addition and subtraction; relate the strategy to a written method and explain the reasoning used. (1-LS1-2) | LS1.A: Structure and Function • All organisms have external parts. Different animals use their body parts in different ways to sense, move, grasp objects, protect themselves, move from place to place, and seek, find, and take in food, water and air. Plants also have different parts (roots, stems, leaves, flowers, buds) that help them survive and grow. (1-LS1-1) LS1.B: Growth and Development of Organisms • Adult plants and animals can have young. In many kinds of animals, parents and the offspring themselves engage in behaviors that help the offspring to survive. (1-LS1-2) LS1.D: Information Processing • Animals have body parts that capture and convey different kinds of information needed for growth and survival. Animals respond to these inputs with behaviors that help them survive. Plants also respond to some external inputs. (1-LS1-1) | Patterns • Patterns in the natural and human designed world can be observed, used to describe phenomena, and used as evidence. (1-LS1-2) Structure and Function • The shape and stability of structures of natural and designed objects are related to their functions. (1-LS1-1) ----- Connections to Engineering, Technology, and Applications of Science Influence of Science, Engineering and Technology on Society and the Natural World • Every human-made product is designed by applying some knowledge of the natural world and is built using materials derived from the natural world. (1-LS1-1) |

How to Make a Simple Electric



★★★★★ 3.8 based on 116 ratings

By Erin Bjornson
Updated on Sep 05, 2013

Energy comes in many forms. Electric energy can be converted into useful work, or mechanical energy, by machines called electric motors. Electric motors work due to electromagnetic interactions: the interaction of current (the flow of electrons) and a magnetic field.

[Download Project](#)

Problem

Find out how to make a simple electric motor.

Materials

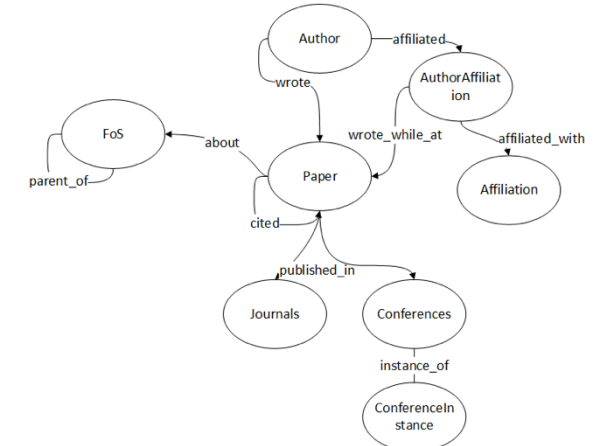
- D battery
- Insulated 22G wire
- 2 large-eyed, long, metal sewing needles (the eyes must be large enough to fit the wire through)
- Modeling clay
- Electrical tape
- Hobby knife
- Small circular magnet
- Thin marker

- Extracted science projects from websites and stored as a structured data

- Extracted all 1188 projects from ScienceBuddies.com

- Extracted DBLP bibliographic database for computer science and MICRO 50 years of publications (~1400)

- All stored in a graph database (~100G) with a structure similar to the Microsoft Academic Graph

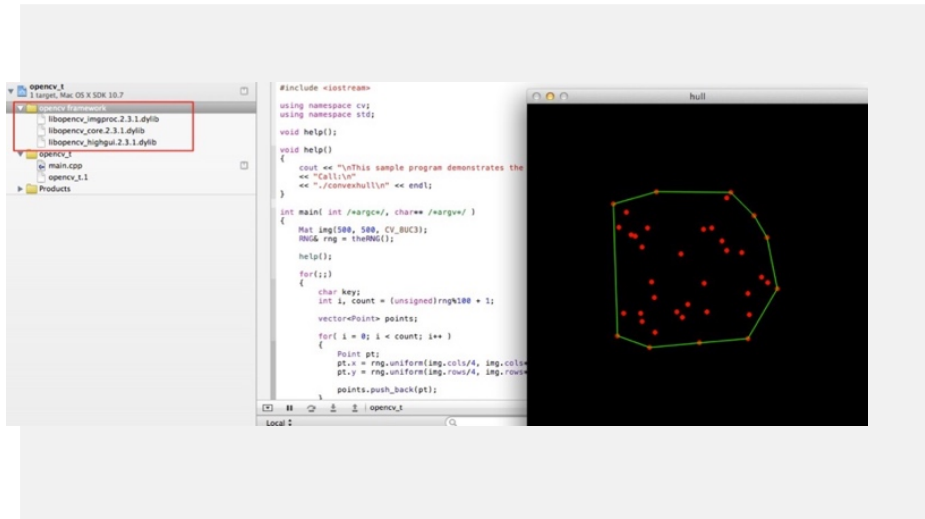


Selected center progress highlights

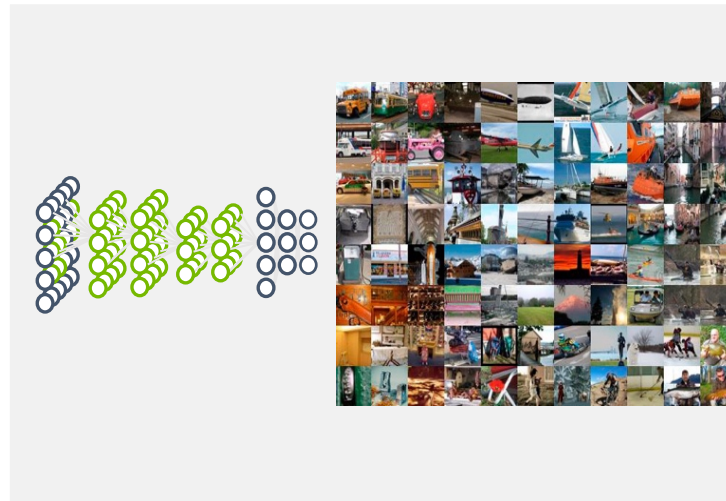
- Curated datasets
- **The CarML system for model development and deployment**
- Workload acceleration
- The Erudite NMA system

Deep Learning Revolution

- a humble beginning in 2010

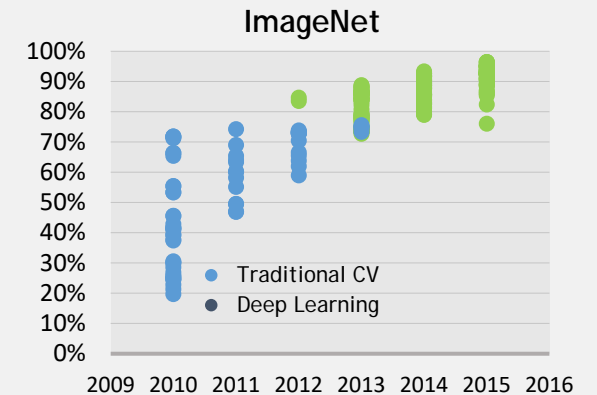


Traditional Computer Vision
Experts + Time



Deep Learning Object Detection
DNN + Data + HPC

2M training images

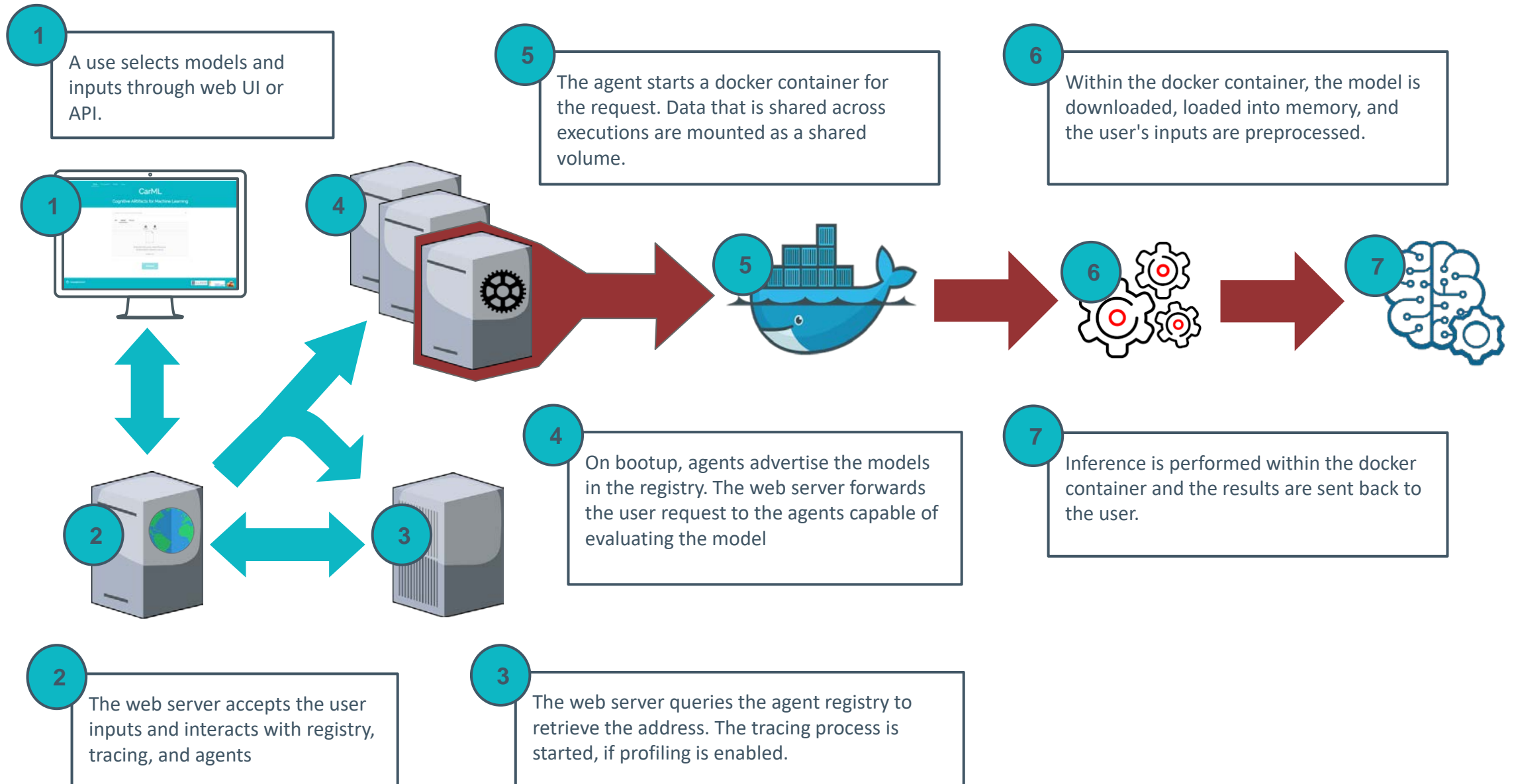


Deep Learning Achieves
"Superhuman" Results

CarML – Cognitive Artifacts for Machine Learning

- CarML.org
- An open source distributed platform to easily deploy and benchmark machine learning frameworks and models **across hardware architectures**, through a common interface.
 - An experimentation platform for ML users
 - A deployment platform for ML developers
 - A benchmarking platform for systems architects

CarML.org as a Web Service



Model Catalog

- Repository contains more than 100 DL models
- Support for Tensorflow, Caffe, Caffe2, and MXNet
 - PyTorch, CNTK, Paddle, ... planned
- Versioned models and frameworks
 - Allows to experiment with custom DL layers

Dataset Catalog

- Repository contains common DL datasets
 - CIFAR 10/100
 - MNIST
 - ImageNet
 -
- Allows one to compare DL models on validation datasets

Machine Catalog

- X86 and Power8 Systems
 - CPU only mode and/or GPU mode
- Planned to have ARM cores and integration with simulators

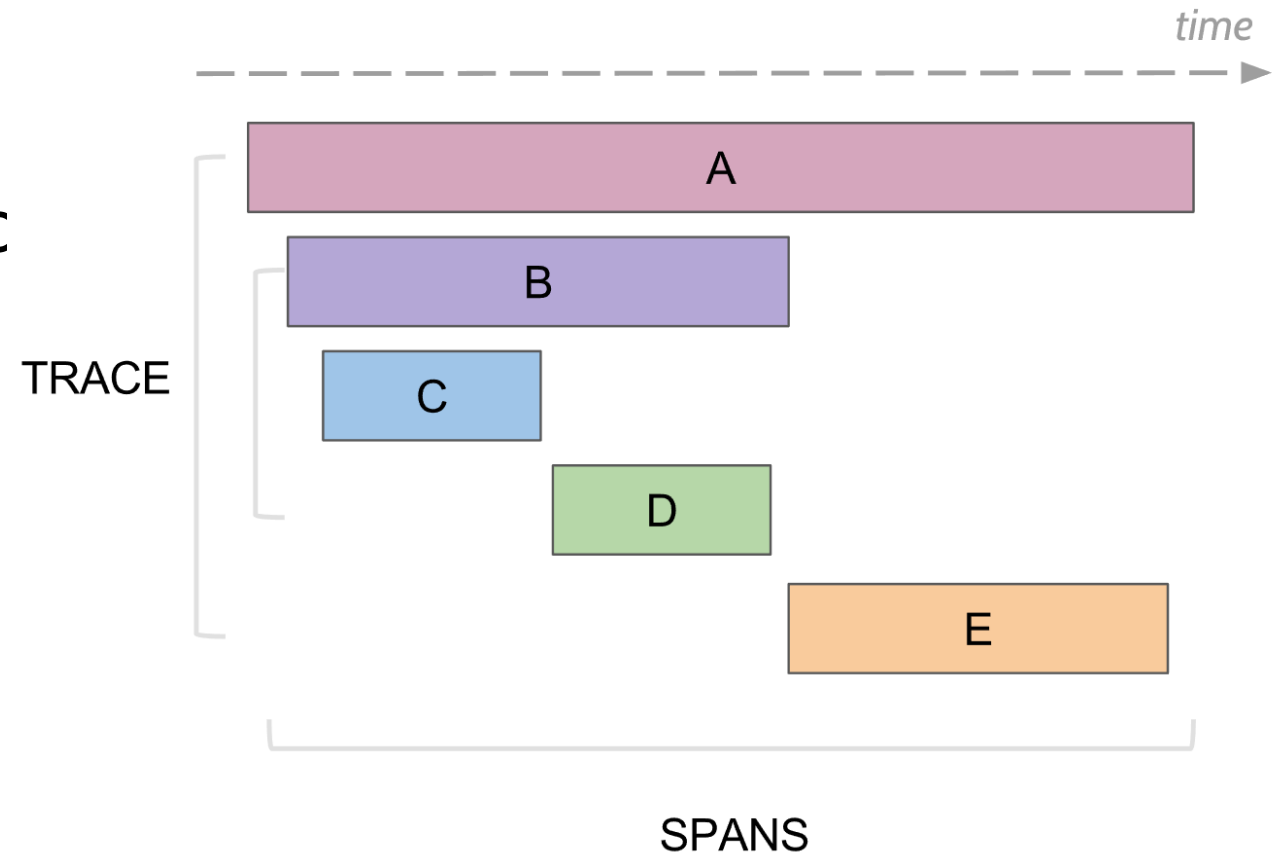
Tracing and Monitoring Options

- Integration with PAPI
- Integration with Perf Events
- Integration with NVIDIA's CUPTI
- Integration with OSX's Instruments

Tracing

Terminology

A **Trace** is a directed acyclic graph (DAG) of **Spans**
Spans can reference one another.



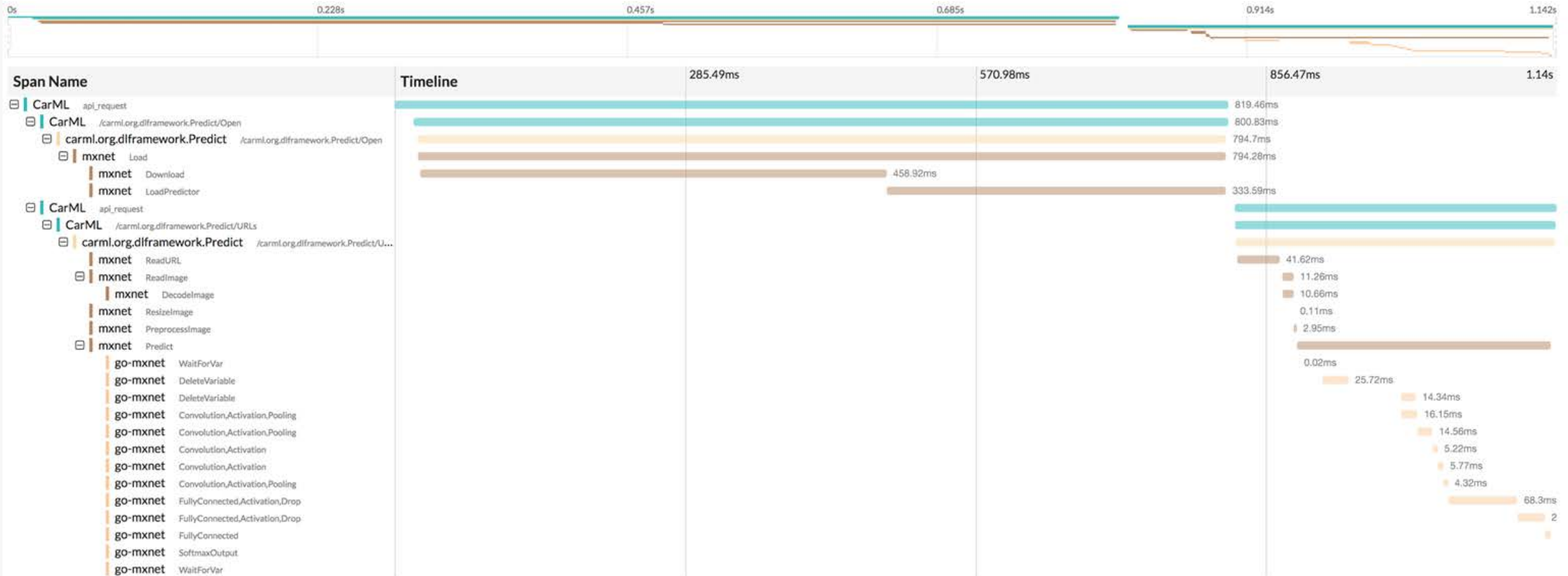
Tracing

CarML: api_request

Trace Start: September 25, 2017 2:21 PM Duration: 1.142s Services: 4 Depth: 6 Total Spans: 28

View Options ▾

Search...



Observers

- Subscribe on StartSpan / EndSpan events
- Capture hardware counters for each event
 - PAPI
 - NVML
 - Perf

CUPTI

- Capture CUDA runtime & driver events
- Integrated with the CarML tracer
 - Implemented in Go
 - Declare CUPTI callback function in Go
 - Pass CUPTI Go handle into C code
 - Events to capture are configurable

CUPTI

```

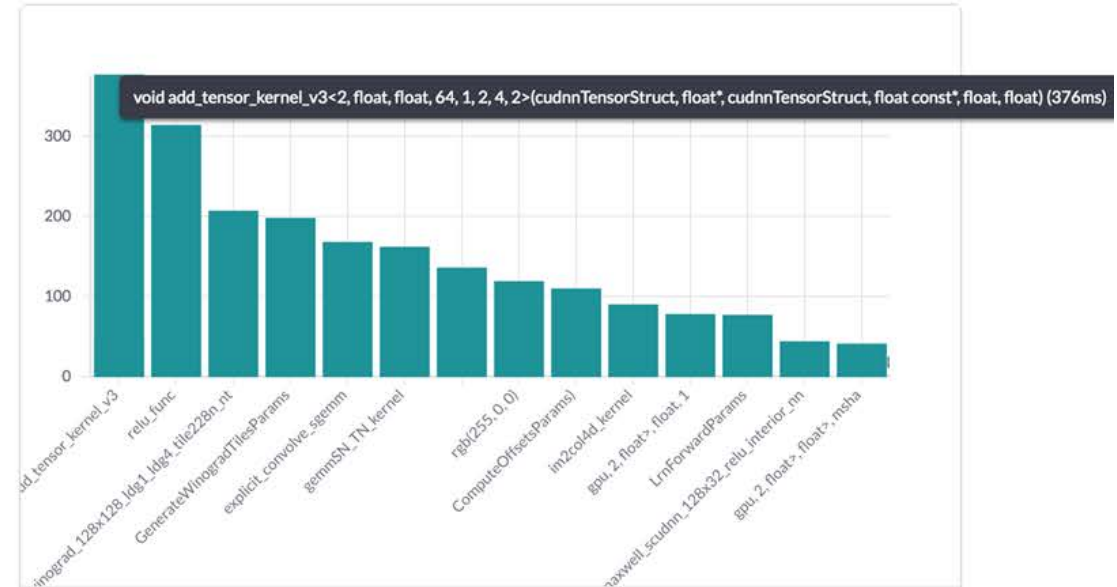
func callback(userData unsafe.Pointer, domain0 C.CUpti_CallbackDomain, cbid0 C.
    handle := (*CUPTI)(unsafe.Pointer(userData))
    if handle == nil {
        log.Debug("expecting a cupti handle, but got nil")
        return
    }
    domain := types.CUpti_CallbackDomain(domain0)
    switch domain {
    case types.CUPTI_CB_DOMAIN_DRIVER_API:
        cbid := types.CUpti_driver_api_trace_cbid(cbid0)
        switch cbid {
        case types.CUPTI_DRIVER_TRACE_CBID_cuLaunchKernel:
            handle.onCULaunchKernel(domain, cbid, cbInfo)
            return
        case types.CUPTI_DRIVER_TRACE_CBID_cuMemcpyHtoD_v2,
            types.CUPTI_DRIVER_TRACE_CBID_cuMemcpyDtoH_v2,
            types.CUPTI_DRIVER_TRACE_CBID_cuMemcpyDtoD_v2,
            types.CUPTI_DRIVER_TRACE_CBID_cuMemcpyHtoDAsync_v2,
            types.CUPTI_DRIVER_TRACE_CBID_cuMemcpyDtoHAsync_v2,
            types.CUPTI_DRIVER_TRACE_CBID_cuMemcpyDtoDAsync_v2:
            handle.onCudaMemcpyDevice(domain, cbid, cbInfo)
            return
        default:
            log.WithField("cbid", cbid.String()).
                WithField("function_name", demangleName(cbInfo.functionName)).
                Debug("skipping runtime call")
            return
        }
    case types.CUPTI_CB_DOMAIN_RUNTIME_API:
        cbid := types.CUPTI_RUNTIME_TRACE_CBID(cbid0)
        switch cbid {
        case types.CUPTI_RUNTIME_TRACE_CBID_cudaDeviceSynchronize_v3020,
            types.CUPTI_RUNTIME_TRACE_CBID_cudaStreamSynchronize_v3020:
            handle.onCudaDeviceSynchronize(domain, cbid, cbInfo)
            return
        case types.CUPTI_RUNTIME_TRACE_CBID_cudaMemcpy_v3020,
            types.CUPTI_RUNTIME_TRACE_CBID_cudaMemcpyAsync_v3020:
            handle.onCudaMemcpy(domain, cbid, cbInfo)
            return
        case types.CUPTI_RUNTIME_TRACE_CBID_cudaLaunch_v3020:
            handle.onCudaLaunch(domain, cbid, cbInfo)
            return
        case types.CUPTI_RUNTIME_TRACE_CBID_cudaThreadSynchronize_v3020:
            handle.onCudaSynchronize(domain, cbid, cbInfo)
            return
        case types.CUPTI_RUNTIME_TRACE_CBID_cudaConfigureCall_v3020:
            handle.onCudaConfigureCall(domain, cbid, cbInfo)
            return
        }
    }
}

```

Current Work

| Filter files | | |
|-----------------------------------|---------|---------------|
| File | Size | Last Modified |
| caffe2_out | | |
| mxnet_out | | |
| BVLC-AlexNet-v1.0 | | |
| mxnet_0.11.0_BVLC-AlexNet_1.0_1 | 162 kB | 3 hours ago |
| mxnet_0.11.0_BVLC-AlexNet_1.0_2 | 169 kB | 3 hours ago |
| mxnet_0.11.0_BVLC-AlexNet_1.0_4 | 182 kB | 3 hours ago |
| mxnet_0.11.0_BVLC-AlexNet_1.0_8 | 189 kB | 3 hours ago |
| mxnet_0.11.0_BVLC-AlexNet_1.0_16 | 148 kB | 3 hours ago |
| mxnet_0.11.0_BVLC-AlexNet_1.0_32 | 259 kB | 3 hours ago |
| mxnet_0.11.0_BVLC-AlexNet_1.0_64 | 320 kB | 3 hours ago |
| mxnet_0.11.0_BVLC-AlexNet_1.0_128 | 455 kB | 3 hours ago |
| mxnet_0.11.0_BVLC-AlexNet_1.0_256 | 700 kB | 3 hours ago |
| mxnet_0.11.0_BVLC-AlexNet_1.0_512 | 1.16 MB | 3 hours ago |
| BVLC-GoogLeNet-v1.0 | | |
| mxnet_0.11.0_BVLC-GoogLeNet_1.0_1 | 851 kB | 3 hours ago |
| mxnet_0.11.0_BVLC-GoogLeNet_1.0_2 | 860 kB | 3 hours ago |
| mxnet_0.11.0_BVLC-GoogLeNet_1.0_4 | 864 kB | 3 hours ago |
| mxnet_0.11.0_BVLC- | 893 | 3 hours |

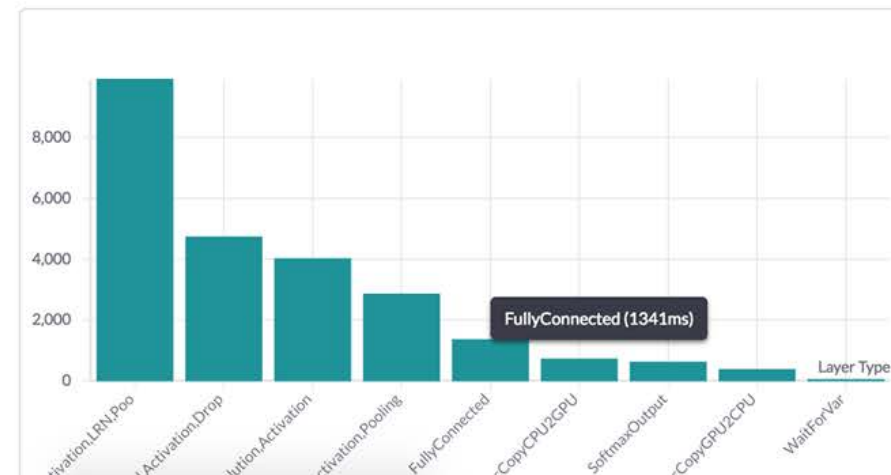
mxnet/BVLC-AlexNet::2



Histogram of CUDA Kernels Executions

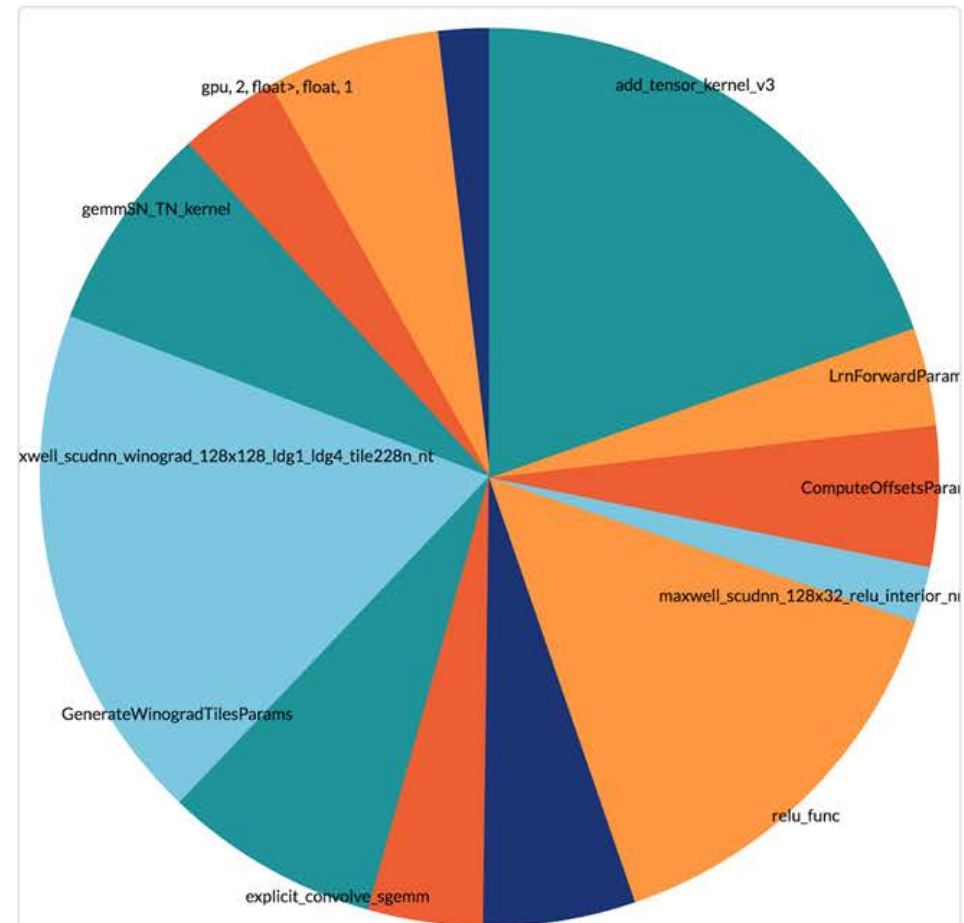
The y-axis is the total time for each CUDA Kernels

Show Pie chart



| Filter files | | |
|-----------------------------------|---------|---------------|
| File | Size | Last Modified |
| caffe2_out | | |
| mxnet_out | | |
| BVLC-AlexNet-v1.0 | | |
| mxnet_0.11.0_BVLC-AlexNet_1.0_1 | 162 kB | 3 hours ago |
| mxnet_0.11.0_BVLC-AlexNet_1.0_2 | 169 kB | 3 hours ago |
| mxnet_0.11.0_BVLC-AlexNet_1.0_4 | 182 kB | 3 hours ago |
| mxnet_0.11.0_BVLC-AlexNet_1.0_8 | 189 kB | 3 hours ago |
| mxnet_0.11.0_BVLC-AlexNet_1.0_16 | 148 kB | 3 hours ago |
| mxnet_0.11.0_BVLC-AlexNet_1.0_32 | 259 kB | 3 hours ago |
| mxnet_0.11.0_BVLC-AlexNet_1.0_64 | 320 kB | 3 hours ago |
| mxnet_0.11.0_BVLC-AlexNet_1.0_128 | 455 kB | 3 hours ago |
| mxnet_0.11.0_BVLC-AlexNet_1.0_256 | 700 kB | 3 hours ago |
| mxnet_0.11.0_BVLC-AlexNet_1.0_512 | 1.16 MB | 3 hours ago |
| BVLC-GoogLeNet-v1.0 | | |
| mxnet_0.11.0_BVLC-GoogLeNet_1.0_1 | 851 kB | 3 hours ago |
| mxnet_0.11.0_BVLC-GoogLeNet_1.0_2 | 860 kB | 3 hours ago |
| mxnet_0.11.0_BVLC-GoogLeNet_1.0_4 | 864 kB | 3 hours ago |
| mxnet_0.11.0_BVLC- | 893 | 3 hours |

mxnet/BVLC-AlexNet::2



PieChart of CUDA Kernels Executions

The area is the total time taken by each CUDA Kernels

☒ Show Pie chart

Model Accuracy on different machines (CPU)

| | mxnet-m | mxnet-l | caffe-m | caffe-l | caffe2-m | caffe2-l |
|-----------------|-----------------------|---------|---------|---------|----------|----------|
| BVLC-AlexNet | 0.4268 | 0.6764 | 0.4268 | 0.6764 | 0.4268 | 0.6764 |
| BVLC-GoogLeNet | 0.9984 | 0.9991 | 0.9984 | 0.9991 | 0.9968 | 0.9991 |
| SqueezeNet-v1.0 | 0.8834 | 0.8501 | 0.7999 | 0.7999 | 0.7999 | 0.9484 |
| SqueezeNet-v1.1 | Shower Cap(0.2874) | 0.6929 | 0.9645 | 0.9645 | 0.9645 | 0.9108 |

m: minsky

l- mac

img: cheeseburger

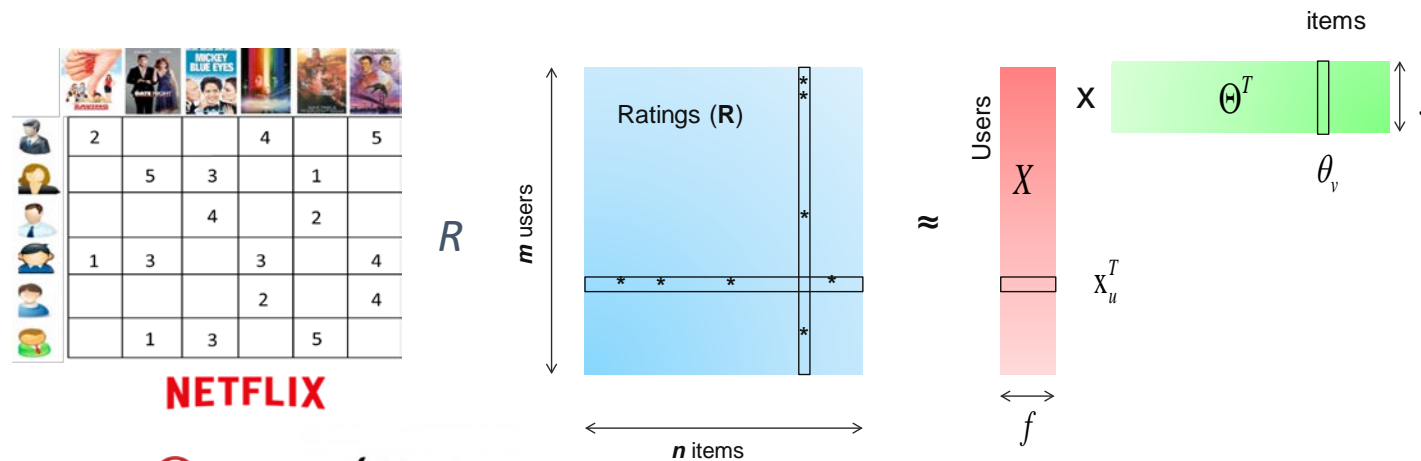
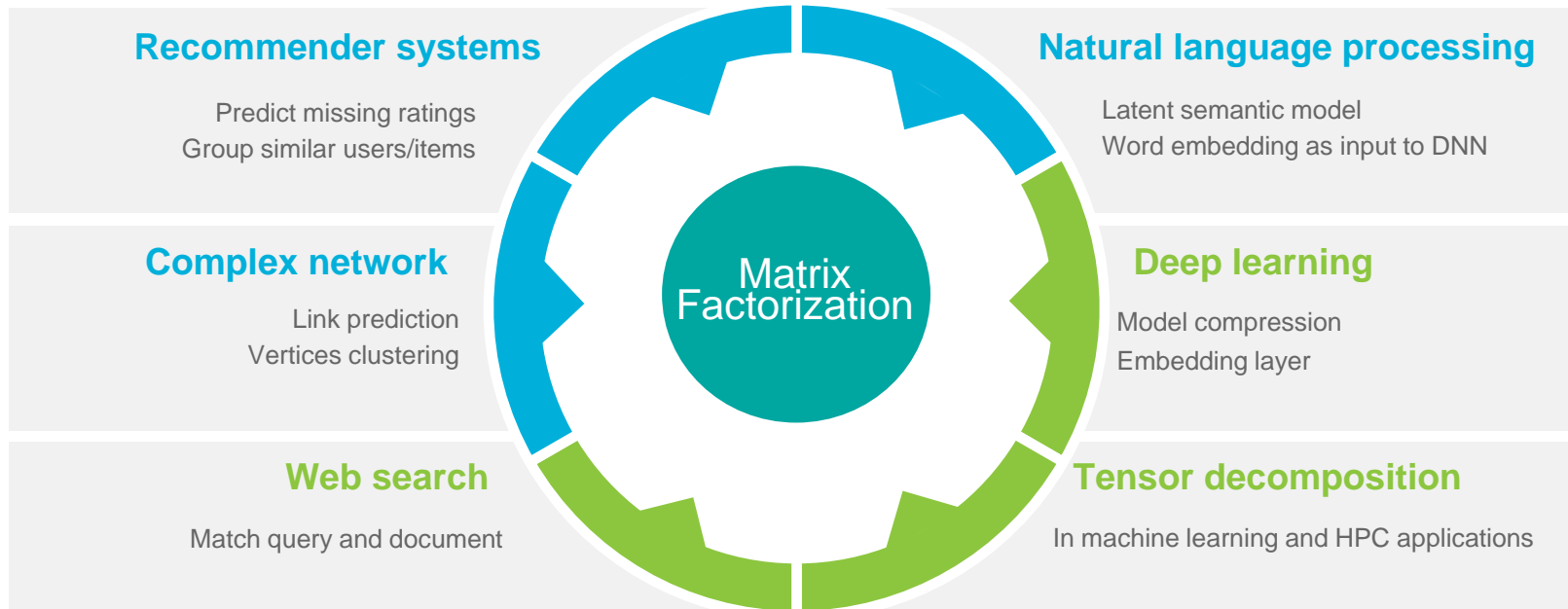
Selected center progress highlights

- Curated datasets
- The CarML System for Model Development and Deployment
- **Workload acceleration**
- The Erudite NMA system

Workload acceleration research at C3SR

- Focus on impactful cognitive workloads for acceleration
 - [Matrix factorization on GPU](#)
 - Long-term Recurrent Convolutional Network acceleration
 - ResNet inference acceleration
 - Neuron Machine Translation acceleration
 - DNN inference acceleration
 - Graph analytic acceleration
- In discussion with other CHN centers to collect performance critical cognitive workloads
- Plan to deliver a set of cognitive benchmarks optimized for OpenPOWER

Matrix factorization: one of key workloads



cuMF acceleration

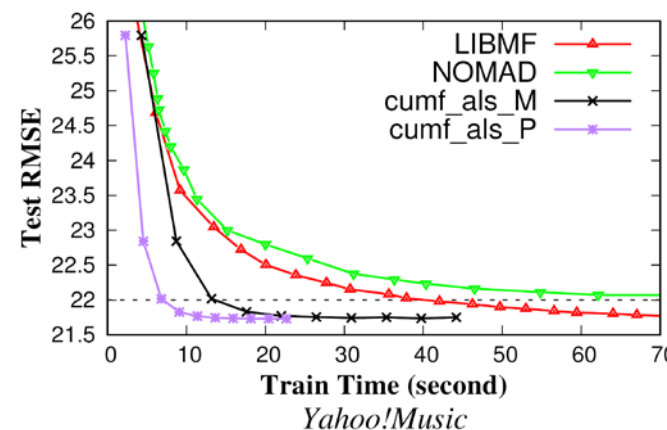
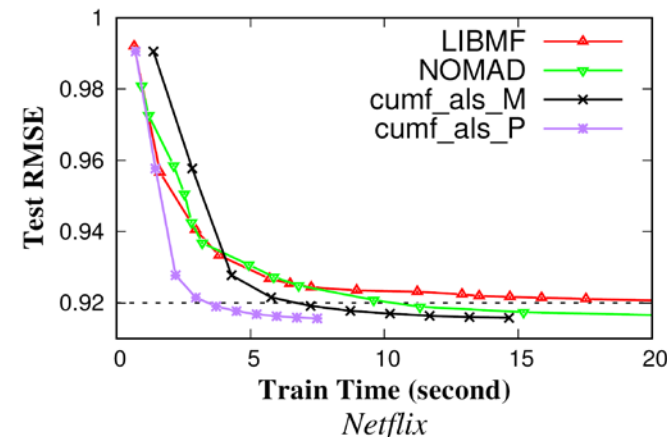
- cuMF formulation: factorize matrix R into

$$R \approx X \cdot \Theta^T$$

- while minimizing the empirical lost

$$J = \sum_{u,v} (r_{uv} - \mathbf{x}_u^T \boldsymbol{\theta}_v)^2 + \lambda (\sum_u n_{x_u} \|\mathbf{x}_u\|^2 + \sum_v n_{\theta_v} \|\boldsymbol{\theta}_v\|^2)$$

- Connect cuMF to Spark MLlib via JNI
- cuMF_ALS @4 Maxwell (\$2.5/hour)
≈ 10x speedup over SparkALS @50 nodes
≈ 1% of SparkALS's cost (\$0.53/hour/node)
- Open source @ <http://github.com/cuMF/>
- **Demoed at SC'16 and GTC'16 on Minsky**
- **Presented to Jen-Hsun Huang on Feb 1, 2017**



- cuMF_ALS w/ FP16 on Maxwell and Pascal
- LIBMF: 1 CPU w/ 40 threads
- NOMAD
 - 32 nodes for Netflix and Yahoo
- 2-10x as fast

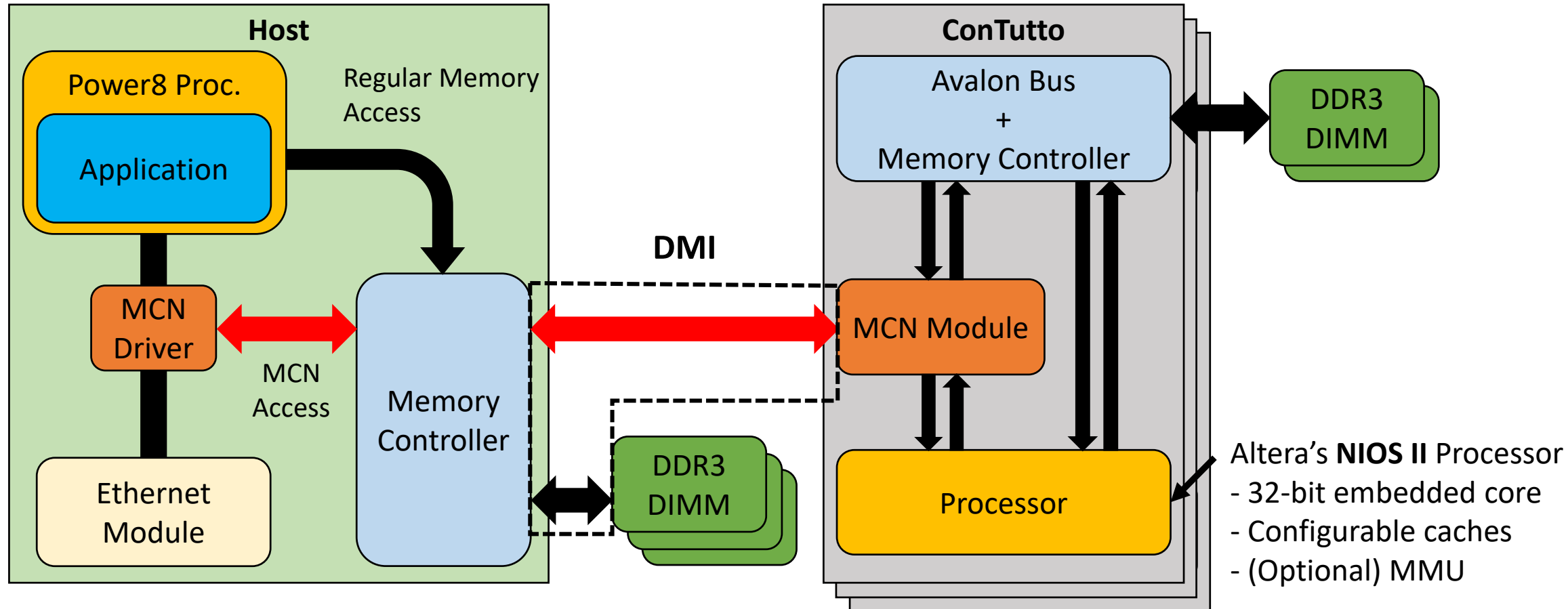
Selected center progress highlights

- Curated datasets
- The CarML System for Model Development and Deployment
- Workload acceleration
- **The Erudite NMA system**

Key Erudite Features

- Persistent objects for main stream languages (C++, Java, Python, etc.)
- Storage-Class Memory (Flash RAM).
- Near Memory Acceleration and memory-channel networking
- API for collaborative CPU/GPU/NMA execution

High-level Diagram of Current MCN Implementation



- Host requires a new kernel driver to transform TCP/IP packet to memory access and vice versa.

Summary

- Creative experiential learning advisor (CELA) as a grand challenge use case for cognitive capabilities
- Cognitive application builder (CAB) to make the underlying heterogeneous infrastructure easy to consume for cognitive application developers
- Cognitive systems innovations (Erudite) for workload acceleration, including Near Memory Acceleration (NMA)