



Accelerating Scientific Applications with SambaNova Reconfigurable Architecture

Vijay Tatkar

Director, SambaNova Systems

BoF at Supercomputing 2024

20 November 2024, Atlanta, GA





Agenda

1. SambaNova: Who We Are
2. Hardware Architecture and Compilation flow
3. Project Redwood: C++ SDK

SambaNova: Who We Are



Who We Are

Snapshot

- Founded in 2017 by industry luminaries and originated at Stanford University
- Fully integrated generative AI platform, from 4th generation hardware to pre-trained models
- \$1B+ funding raised

Founded by pioneers in AI



Lip-Bu Tan
Executive Chairman



Rodrigo Liang
Co-founder & CEO



Kunle Olukotun
*Co-founder &
Chief Technologist &
Stanford Professor*



Christopher Ré
*Co-founder & Stanford
Professor*

Sophisticated, long-term investors

BlackRock
Capital Investment Corporation™

SoftBank
Investment Advisors

TEMASEK



SAMSUNG
CATALYST
FUND





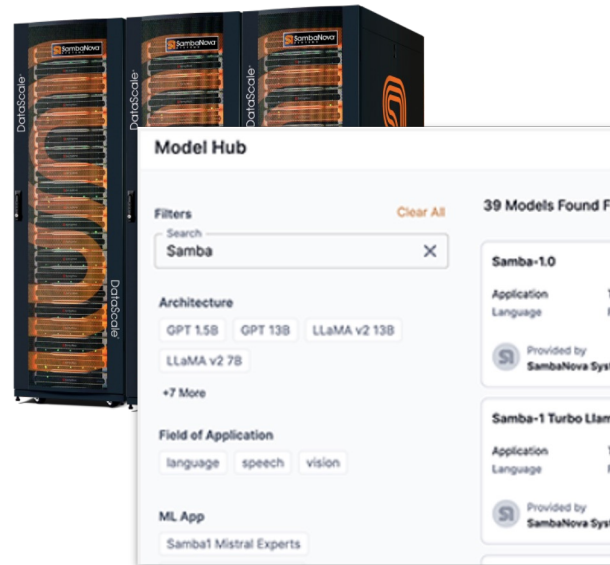
SambaNova Products

SambaNova Cloud



API service for inference at record-breaking speeds ([link](#))

SambaNova Suite



Secure, on-premises AI platform for training and inference ([link](#))

SambaNova DataScale



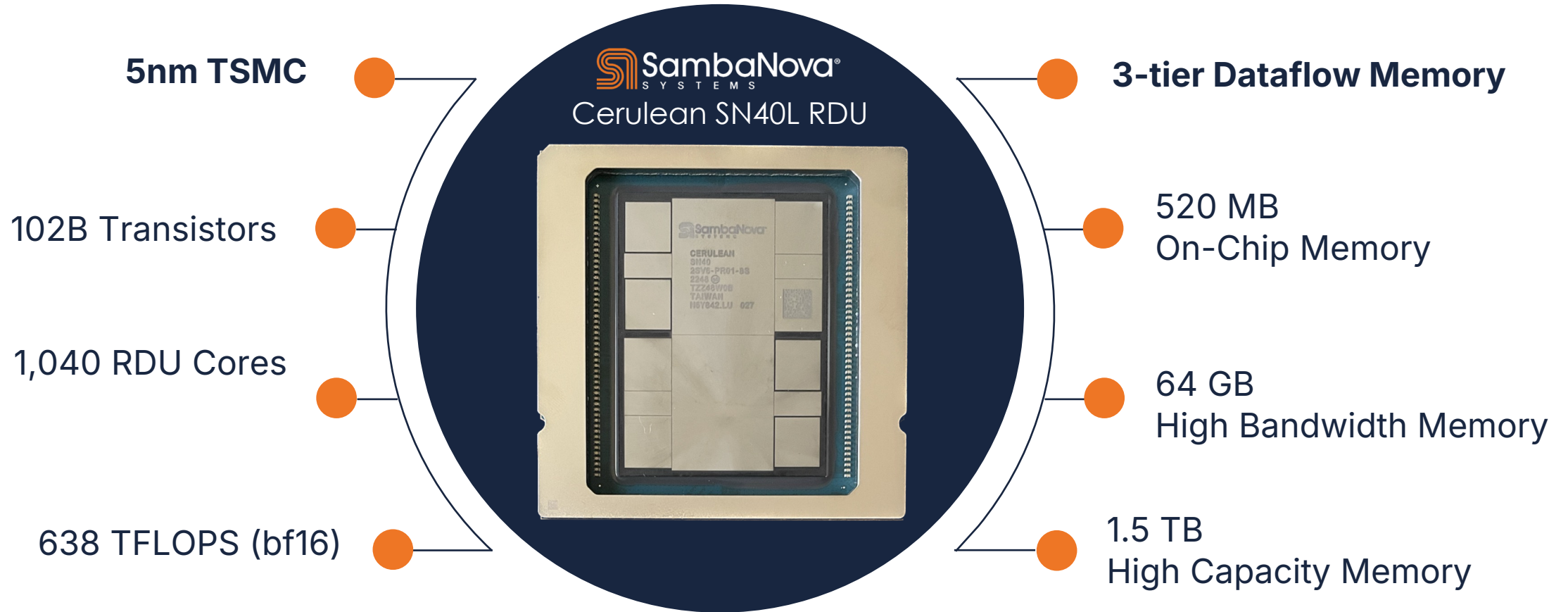
Fully integrated hardware-software AI system ([link](#))

Hardware and Software Architecture



SN40L: SambaNova's 4th Gen AI Chip

"Cerulean" Architecture-based Reconfigurable Dataflow Unit

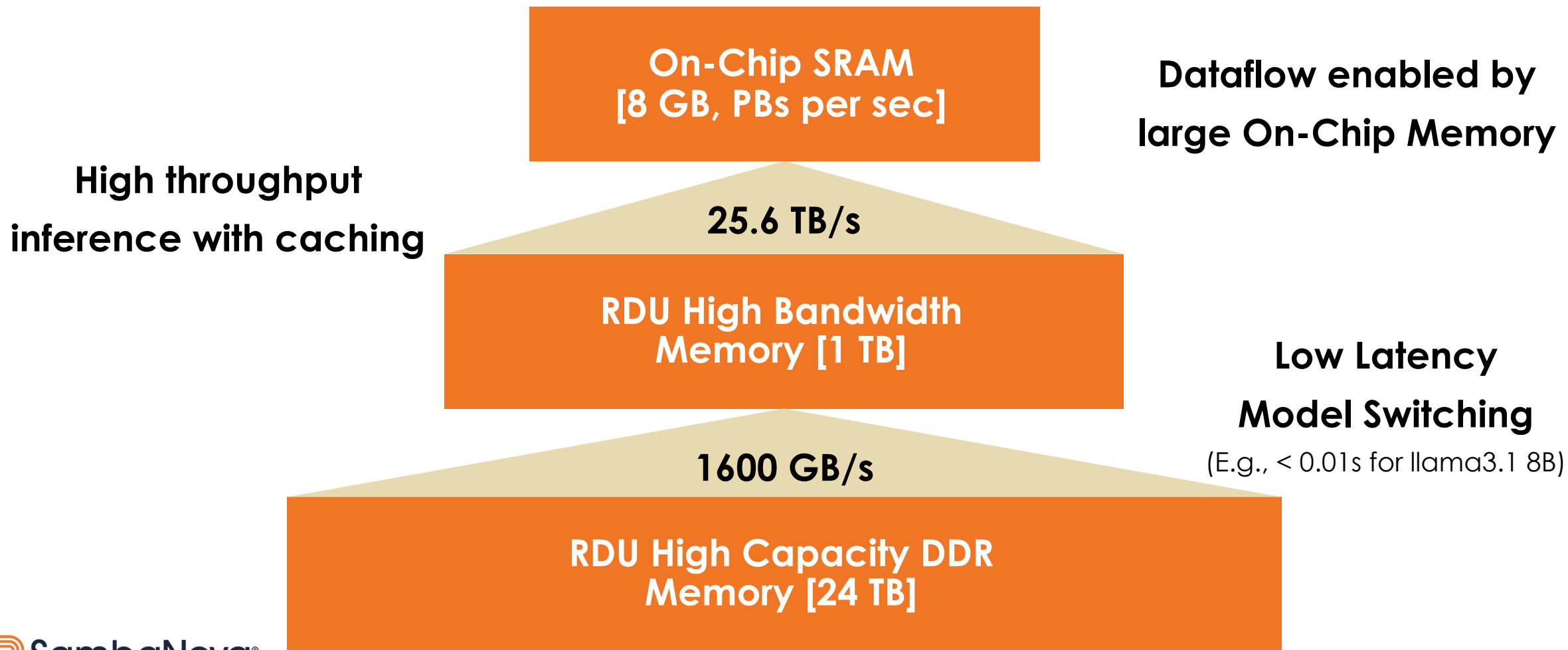


Generative AI Training and Inference



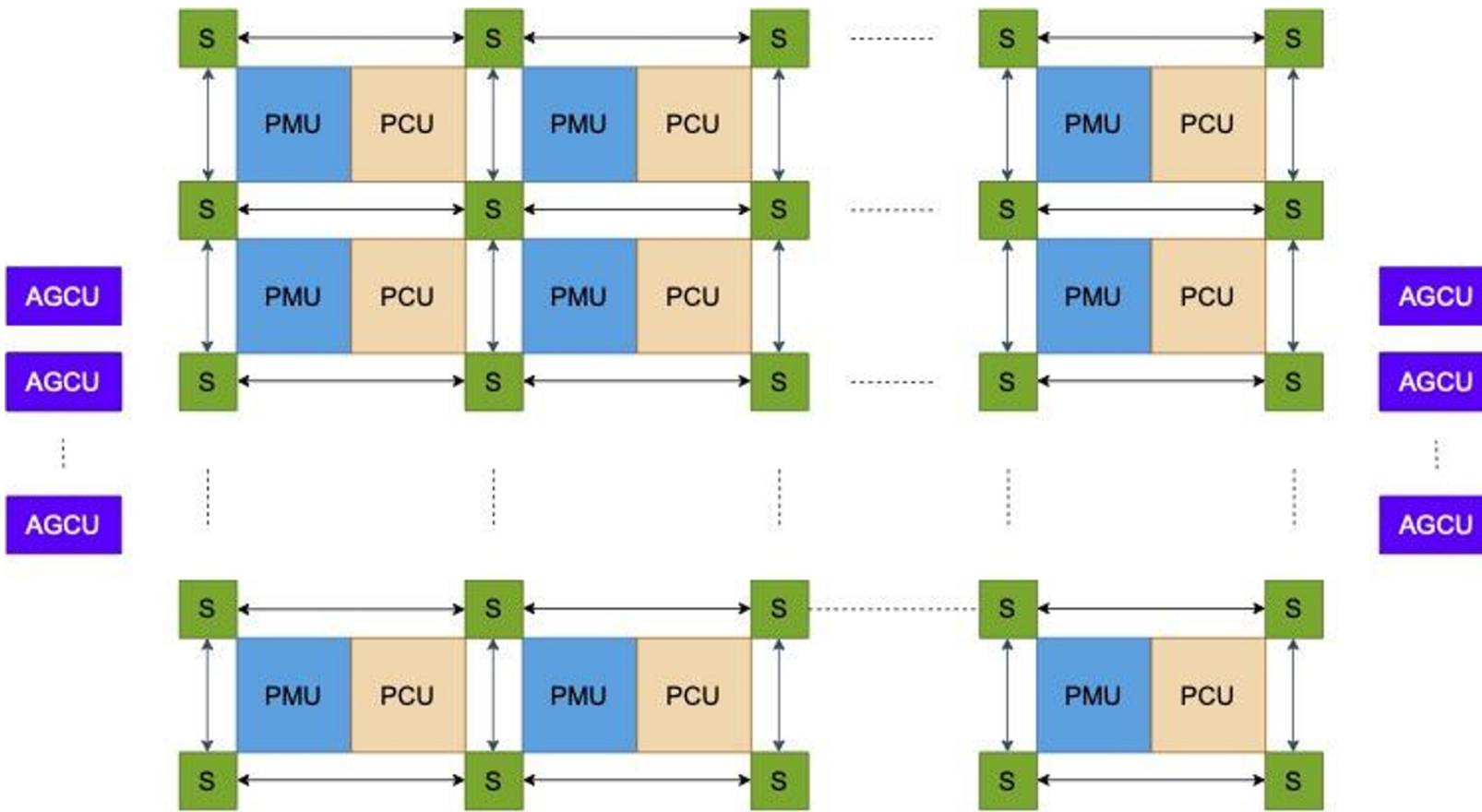
SN40L: Three Tier Memory Architecture

3-tier Memory System with SRAM, HBM, and DDR





SN40L: Tile Architecture



1040 PCUs and PMUs

PCU: Compute unit

PMU: Memory unit

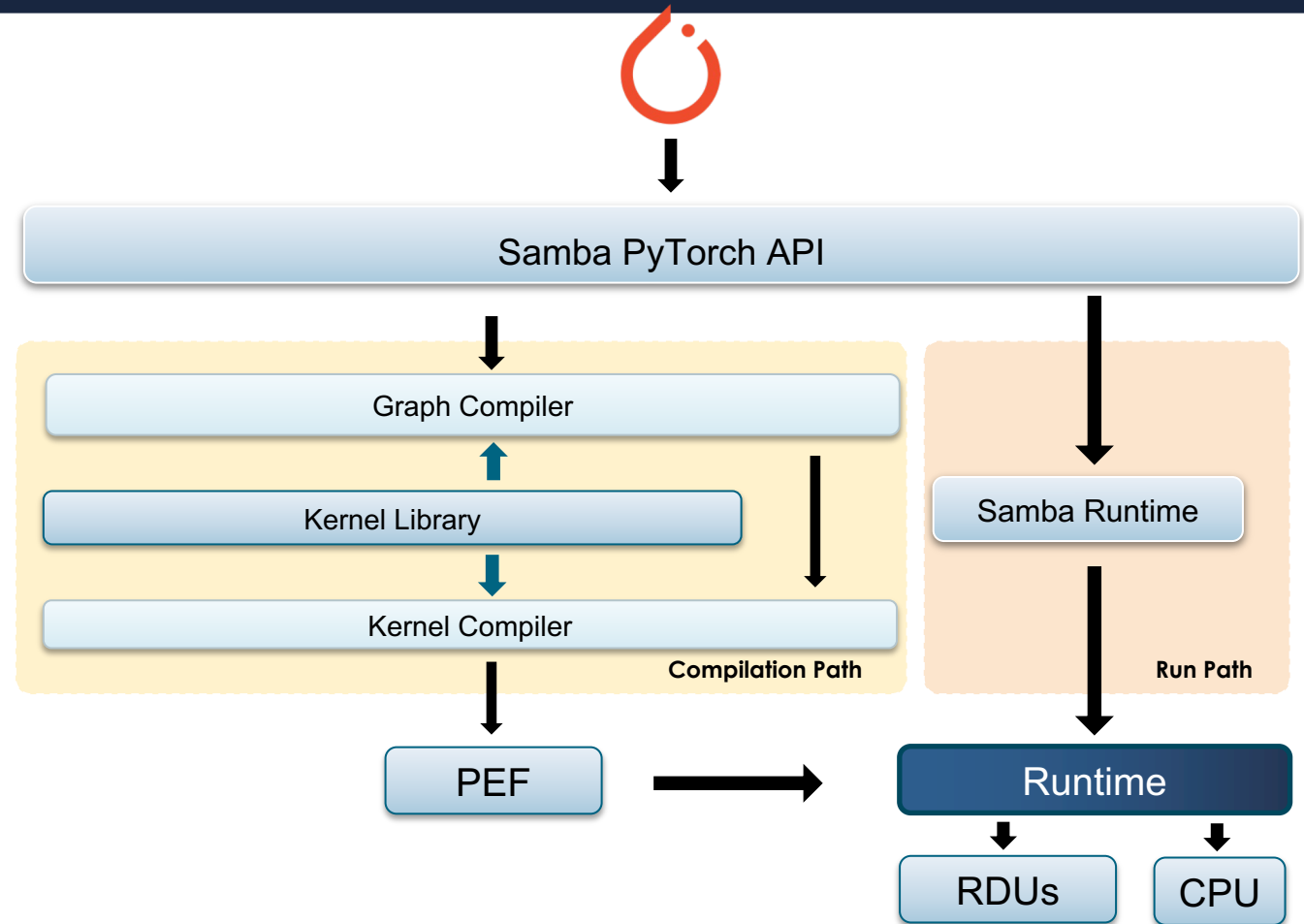
S: Mesh switches

AGCU: Portal to off-chip memory and IO



Samba Compilation Flow

- **Samba**
 - + SambaNova PyTorch compilation & run APIs
- **Graph compiler**
 - + High-level ML graph transformation & optimizations
- **Kernel compiler**
 - + Low-level RDU operator kernel transformation & optimizations
- **Kernel library**
 - + RDU operator implementations



Project Redwood

C++ SDK



Redwood: C++ SDK for the RDU

- What is Project Redwood:
 - Tensor-oriented kernel definition language and RDU scheduling SDK, embedded in standard C++
- Redwood lets users
 - + **Specify tensor functions** in standard C++
 - + **Compile them from** C++ API
 - + **Run them from** C++ API
 - + **Tune** them from C++, aided by SambaTune
 - + **Debug** them from C++ through emulation, watchpoints and alerts



Redwood: Design Objectives and Status

Design Objectives

- Enable expert developers to exploit the capabilities of RDUs
- For new innovation vs. porting
- Example use cases
 - + Convert compute heavy inner loops of existing C++ programs as tensor for RDU offload
 - + Develop high-performance ML operators

Status

- Ramping internal use
- Early preview with select customers
- Feedback collection will inform design choices for needed kernels
- Public release coming soon



Redwood: Goals and Programming Model

Goals for Redwood library

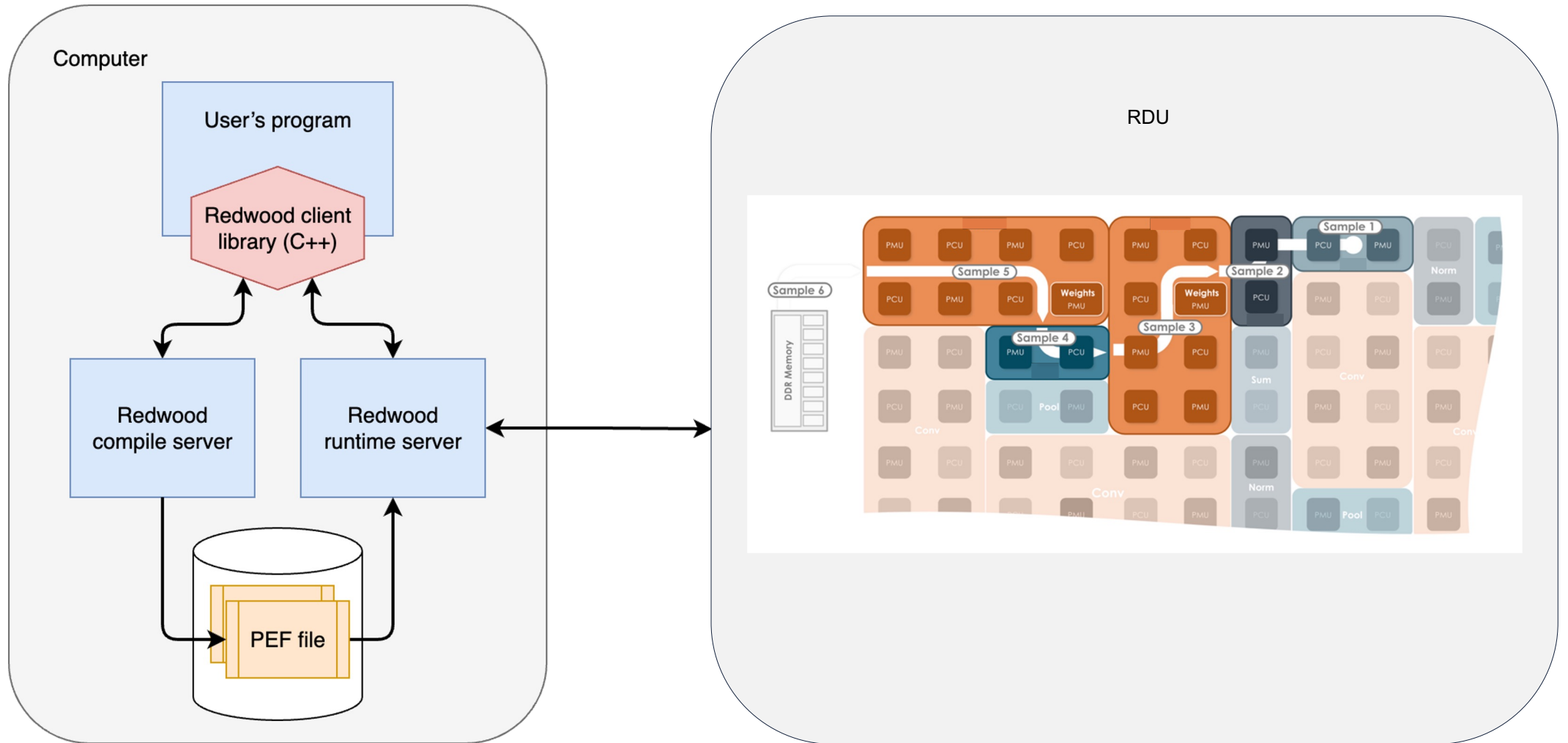
- Simplicity of numpy arrays
- Predictable performance characteristics of Fortran
- Leverage compiler to exploit parallel patterns (map, filter, reduce) to specify loop-like constructs
- Optimization directives (fusion, tiling, parallelization) for power users
- Composability, reusability and zero-cost abstractions

Redwood Programming Model

- Redwood tensor is “abstract”. SDK executes kernels symbolically
- Redwood array is concrete and used for data motion
- Tensors expose their statically known shapes; allows for implementation choice based on extent of dimension
- Any function that manipulates `redw::Tensor` can be a kernel; function calls have no overhead in binary



Redwood: System Components





SambaNova Cloud

Putting large scale applications together



	SambaNova	GPU
Llama 3.2 1B 16-bit	2477	304
Llama 3.1 8B 16-bit	1066	93
Llama 3.1 70B 16-bit	460	32
Llama 3.1 405B 16-bit	200	14

THANK YOU!

vijay.tatkar@sambanova.ai

Join us

Booth #2309

sambanova.ai/sc24

Meet the Experts and Happy Hours at Booth #2309

Sunday, November 17

- 11:30 a.m. - 12:00 noon

Tutorial: GenAI Training and Inference at Scale

Tuesday, November 19

- 1:00 p.m. - 3:30 p.m.

SambaNova Customer Experts

- 4:00 p.m. - 5:00 p.m.

SambaNova Experts Happy Hour

- 5:00 p.m. - 6:00 p.m.

SambaNova Partner Experts Happy Hour

Wednesday, November 20

- 1:00 p.m. - 2:30 p.m.

Meet SambaNova Customer Experts

Birds of a Feather

Wednesday, November 20

- 5:15 p.m. - 6:45 p.m.

*Democratizing AI Accelerators for HPC Applications:
Challenges, Success, and Support*

- 5:15 p.m. - 6:45 p.m.

*The National Artificial Intelligence Research Resource (NAIRR)
Pilot User Experience BoF*



Try It Today
cloud.sambanova.ai



Questions? Join the Community

The screenshot shows the SambaNova community website. At the top left is the SambaNova logo. A search bar is located at the top center. On the top right, there are icons for a document, a message, a notification bell with a '6' badge, and a user profile picture. The main heading is "Accelerate Your AI Journey with SambaNova!". Below this is a sub-heading: "Unleash the power of AI with SambaNova. Get Started with our documentation and further accelerate your journey with our AI Starter Kits. Network with our amazing developers who are all building the future of AI apps today." There are four orange buttons: "Welcome to the Community", "Documentation", "Discussion", and "Showcase". On the left side, there is a navigation menu with sections: "Mentions", "Bookmarks", "Messages", "Admin", "Categories", "Welcome", "Events", "SambaNova Documentation" (with sub-links for "SambaNova Cloud, Starter Kits", "Fast API Docs", and "Documentation Staging"), and "SambaNova Devs" (with sub-links for "Introductions", "Discussion", and "Showcase"). The main content area features a "Showcase" section with a header "We would love to see what you have built! Showcase it with the broader community." Below this are two posts: one from "FAB (Slack Bot)" dated "3 days ago" with 5 likes and 1 comment, and another from "LE" dated "Sep 11" titled "GoogleSheets and GoogleDocs integration" with 8 likes and 1 comment. On the right side, there is a user profile card for "Hi, vasanth.mohan!" with an "Edit" button, and a "Top Topics" section with a "Quarterly" filter and a list item: "1. A desktop robot integrated with SambaNova's Fast API" dated "11d" with 23 likes and 4 comments.

Community.SambaNova.ai