### **Polynesia:**

Enabling High-Performance and Energy-Efficient Hybrid Transactional/Analytical Databases with Hardware/Software Co-Design

### Amirali Boroumand Geraldo F. Oliveira

Saugata Ghose Onur Mutlu

ICDE 2022



# **Executive Summary**

- <u>Context</u>: Many applications need to perform real-time data analysis using an <u>Hybrid Transactional/Analytical Processing (HTAP)</u> system
  - An ideal HTAP system should have three properties:
    - (1) data freshness and consistency, (2) workload-specific optimization,
    - (3) performance isolation
- <u>Problem</u>: Prior works cannot achieve all properties of an ideal HTAP system
- <u>Key Idea</u>: Divide the system into transactional and analytical processing islands
  - Enables workload-specific optimizations and performance isolation
- <u>Key Mechanism</u>: Polynesia, a novel hardware/software cooperative design for in-memory HTAP databases
  - Implements custom algorithms and hardware to reduce the costs of data freshness and consistency
  - Exploits **PIM** for analytical processing to alleviate data movement
- <u>Key Results</u>: Polynesia outperforms three state-of-the-art HTAP systems
  - Average transactional/analytical throughput improvements of 1.7x/3.7x
  - 48% reduction on energy consumption

#### SAFARI

1	Introduction	
2	Limitations of HTAP Systems	
3	Polynesia: Overview	
4	<b>Update Propagation Mechanism</b>	
5	<b>Consistency Mechanism</b>	
6	Analytical Engine	
7	Evaluation	
8	Conclusion	
	Motivation Polynesia Update Propagation Consistency Mechanism Analytical Engine Evaluation Co	nclusion

3

•

1	Introduction
2	Limitations of HTAP Systems
3	Polynesia: Overview
4	Update Propagation Mechanism
5	<b>Consistency Mechanism</b>
6	Analytical Engine
7	Evaluation
8	Conclusion
SAFARI Introduction	Motivation Polynesia Update Propagation Consistency Mechanism Analytical Engine Evaluation

Conclusion

•

# **Real-Time Analysis**

An explosive interest in many applications domains to perform data analytics on the most recent version of data (real-time analysis)



### For these applications, it is critical to analyze the transactions in real-time as the data's value diminishes over time

SAFARI Introduction

Motivation

n Polynesia

Update Propagation

on Consistency Mechanism

m Analytical Engine

Evaluation

### **HTAP: Supporting Real-Time Analysis**

Traditionally, new transactions (updates) are propagated to the analytical database using a periodic and costly process



To support real-time analysis: a single hybrid DBMS is used to execute both transactional and analytical workloads

Introduction SAFAR . . .

Motivation . . . . . . . .

**Polynesia** • •

Update Propagation . . . . . .

**Consistency Mechanism** • •

**Analytical Engine** . . .

Evaluation Conclusion b

# **Ideal HTAP System Properties**

### An ideal HTAP system should have three properties:

### **Workload-Specific Optimizations**

Transactional and analytical workloads must benefit from their own specific optimizations

#### 2 **Data Freshness and Consistency Guarantees**

 Guarantee access to the most recent version of data for analytics while ensuring that transactional and analytical workloads have a consistent view of data

### **Performance Isolation**

Latency and throughput of transactional and analytical workloads are the same as if they were run in isolation

### Achieving all three properties at the same time is very challenging

Introduction . . .

Motivation

Polynesia

• •

Update Propagation . . . . . .

Consistency Mechanism Analytical Engine • •

Evaluation

1	Introduction	
2	Limitations of HTAP Systems	
3	Polynesia: Overview	
4	Update Propagation Mechanism	
5	<b>Consistency Mechanism</b>	
6	Analytical Engine	
7	Evaluation	
8	Conclusion	
SAFARI Introduction	Motivation Polynesia Update Propagation Consistency Mechanism Analytical Engine Evaluation Consistency Mechanism	onclusion

8

# State-of-the-Art HTAP Systems

### We study two major types of HTAP systems:



**Single-Instance** 

**Multiple-Instance** 

#### We observe two key problems:

**Data freshness and consistency mechanisms** are costly and cause a drastic reduction in throughput

SAFAR

These systems fail to provide performance isolation because of high resource contention

troduction . . .

Motivation . . . . . . . .

**Polynesia** 

• •

**Update Propagation** . . . . . .

Consistency Mechanism • •

**Analytical Engine** . . .

Evaluation Conclusion 9

# State-of-the-Art HTAP Systems

### We study two major types of HTAP systems:



Single-Instance

**Multiple-Instance** 

#### We observe two key problems:

Data freshness and consistency mechanisms are costly and cause a drastic reduction in throughput

2



SAFARI Introduction

Motivation

Polynesia Update

Update Propagation

Consistency Mechanism

Analytical Engine

Evaluation

# Single-Instance: Data Consistency

Since both analytics and transactions work on the same data concurrently, we need to ensure that the data is consistent

There are two major mechanisms to ensure consistency:



### Drawbacks of Snapshotting and MVCC

We evaluate the throughput loss caused by Snapshotting and MVCC:



Throughput loss comes from <u>memcpy</u> operation: generates a large amount of data movement



Throughput loss comes from long version chains: expensive time-stamp comparison and

a large number of random memory accesses

Consistency Mechanism A

Analytical Engine Evaluation

# State-of-the-Art HTAP Systems

### We study two major types of HTAP systems:



**Multiple-Instance** 

#### We observe two key problems:



Introduction SAFAR . . .

Motivation . . . . . . . .

**Polynesia** 

• •

. . . . . .

Update Propagation

Consistency Mechanism • •

**Analytical Engine** . . .

Evaluation

# **Maintaining Data Freshness**

One of the major challenges in multiple-instance systems is to keep analytical replicas up-to-date

**Transactional queries** 



**Multiple-Instance HTAP System** 

To maintain data freshness (via Update Propagation):

- **Update Gathering and Shipping: gather updates from** transactional threads and ship them to analytical the replica
- 2 Update Application: perform the necessary format conversation and apply those updates to analytical replicas

SAFARI Introductio

Motivation

Polynesia

Update Propagation

Consistency Mechanism

nism Analytical Engine

Evaluation

# **Cost of Update Propagation**

### We evaluate the throughput loss caused by Update Propagation:



Transactional <u>throughput reduces</u> by up to <u>21.2%</u> during the update gathering & shipping process

# Transactional <u>throughput reduces</u> by up to <u>64.2%</u> during the update application process

SAFARI Introduction

Motivation

Polynesia U

Update Propagation

agation Consiste

Consistency Mechanism

Analytical Engine

ine Evaluation

# Problem and Goal

### **Problems:**

4	State-of-the-art HTAP systems do not achieve
	all of the desired HTAP properties

Data freshness and consistency mechanisms are 2 data-intensive and cause a drastic reduction in throughput

These systems fail to provide performance isolation ſ because of high resource contention

### Goal:

### Take advantage of custom algorithm and processing-in-memory (PIM) to address these challenges

SAFARI Introduction . . .

Motivation . . . . . . . . Polynesia • •

Update Propagation . . . . . .

Consistency Mechanism Analytical Engine • •

. . .

Evaluation . . . . . . . .

1	Introduction	
2	Limitations of HTAP Systems	
3	Polynesia: Overview	
4	Update Propagation Mechanism	
5	<b>Consistency Mechanism</b>	
6	Analytical Engine	
7	Evaluation	
8	Conclusion	
SAFARI Introduction	Motivation Polynesia Update Propagation Consistency Mechanism Analytical Engine Evaluation Co	onclusion •

7

# Polynesia

SAFARI Introduction

. . .

Key idea: partition computing resources into two types of isolated and specialized processing islands

Isolating transactional islands from analytical islands allows us to:

- Apply workload-specific optimizations to each island
- **Avoid high resource contention**

Motivation

Polynesia

• •

- Design efficient data freshness and consistency mechanisms without incurring high data movement costs
  - Leverage processing-in-memory (PIM) to reduce data movement

Consistency Mechanism Analytical Engine

. . .

• •

Evaluation

Conclusion

**PIM** mitigates data movement overheads by • placing computation units nearby or inside memory

Update Propagation

. . . . . .

# **Polynesia: High-Level Overview**

Each island includes (1) a replica of data, (2) an optimized execution engine, and (3) a set of hardware resources



Introduction SAFARI . . .

Motivation

**Polynesia** • •

**Update Propagation** . . . . . .

**Consistency Mechanism** • •

. . .

**Analytical Engine** 

Evaluation

1	Introduction	
2	Limitations of HTAP Systems	
3	Polynesia: Overview	
4	<b>Update Propagation Mechanism</b>	
5	Consistency Mechanism	
6	Analytical Engine	
7	Evaluation	
8	Conclusion	
SAFARI Introduction	Motivation Polynesia Update Propagation Consistency Mechanism Analytical Engine Evaluation Conclusion	20

# **Maintaining Data Freshness**

One of the major challenges in multiple-instance systems is to keep analytical replicas up-to-date

**Transactional queries** 



**Multiple-Instance HTAP System** 

To maintain data freshness (via Update Propagation):

- Update Gathering and Shipping: gather updates from transactional threads and ship them to analytical the replica
- **Update Application:** perform the necessary format conversation and apply those updates to analytical replicas

SAFARI

Polynesia . .

Update Propagation . . . . . .

Consistency Mechanism • •

Analytical Engine

Evaluation

### **Update Gathering & Shipping: Algorithm**

Update gathering & shipping algorithm has three major stages:



2<sup>nd</sup> and 3<sup>rd</sup> stages generate a <u>large amount of data movement</u> and account for <u>87.2%</u> of our algorithm's execution time

SAFARI Introduction

Motivation

n Polynesia

Update Propagation

Consistency Mechanism

m Analytical Engine

Evaluation

### **Update Gathering & Shipping: Hardware**

### To avoid these **bottlenecks**, we design a new hardware accelerator, called update gathering & shipping unit



### **Update Propagation: Update Application**

Goal: perform the necessary format conversation and apply transactional updates to analytical replicas



A simple tuple update in row-wise layout leads to multiple random accesses in column-wise layout

2 Updates change encoded value in the dictionary → (I) Need to reconstruct the dictionary, and (2) recompress the column SAFARI Introduction Motivation Polynesia Update Propagation Consistency Mechanism Analytical Engine Evaluation Conclusion 24

# **Update Application: Algorithm**

### We design our update application algorithm to be aware of PIM logic characteristics and constraints



# **Update Application: Hardware**

We design a hardware implementation of our algorithm, and add it to each in-memory analytical island



SAFARI Introduction . . .

Motivation

Polynesia • •

Update Propagation . . . . . .

**Consistency Mechanism** • •

**Analytical Engine** 

Evaluation

Conclusion

26

1	Introduction
2	Limitations of HTAP Systems
3	Polynesia: Overview
4	Update Propagation Mechanism
5	Consistency Mechanism
6	Analytical Engine
7	Evaluation
8	Conclusion
SAFARI Introduction	Motivation Polynesia Update Propagation Consistency Mechanism Analytical Engine Evaluation Conclusion 27

### **Consistency Mechanism: Algorithm**

### For each column, there is a chain of snapshots where each chain entry corresponds to a version of the column



### Polynesia creates a new snapshot only if (1) any of the columns are dirty, and (2) no current snapshot exists for the same column

Introduction SAFAR . . .

Motivation . . . . . . . .

Polynesia

• •

Update Propagation . . . . . .

**Consistency Mechanism** • •

. . .

Analytical Engine

Evaluation Conclusion

28

### **Consistency Mechanism: Hardware**

Our algorithm success at satisfying performance isolation relies on how fast we can do memcpy to minimize snapshotting latency



come back from memory out of order. Allows to immediately initiate a write after a read is complete

• •

Introduction SAFAR . . .

Motivation

Polynesia Update Propagation . . . . . .

**Consistency Mechanism** • •

**Analytical Engine** 

Evaluation

1

1	Introduction	
2	Limitations of HTAP Systems	
3	Polynesia: Overview	
4	<b>Update Propagation Mechanism</b>	
5	<b>Consistency Mechanism</b>	
6	Analytical Engine	
7	Evaluation	
8	Conclusion	
SAFARI Introduction	Motivation Polynesia Update Propagation Consistency Mechanism Analytical Engine Evaluation Consistency Mechanism	Conclusion 30



SAFARI Introduction

Motivation

• • •

Update Propagation

Consistency Mechanism

Analytical Engine

Evaluation

### **Analytical Engine: Data Placement**

### Problem: how to partition analytical data across vaults of the 3D-stacked memory



Introduction . . .

SAFAR

Motivation

**Update Propagation** . . . . . .

• •

Consistency Mechanism • •

**Analytical Engine** 

Evaluation

### **Analytical Engine: Query Execution**

Other details in the paper:

Task scheduling policy

We design a pull-based task assignment strategy, where PIM threads cooperatively pull tasks from the task queue at runtime

### How each physical operator is executed

We employ the top-down Volcano (Iterator) execution model to execute physical operations (e.g., scan, filter, join) while respecting operator's dependencies



Motivation . . . . . . . .

**Polynesia** • •

Update Propagation . . . . . .

Consistency Mechanism Analytical Engine • •

. . .

Evaluation

# **Analytical Engine: Query Execution**

**Other details in the paper:** 

**Polynesia: Enabling High-Performance and Energy-Efficient Hybrid Transactional/Analytical Databases** with Hardware/Software Co-Design

Amirali Boroumand<sup>†</sup> <sup>†</sup>Google

Saugata Ghose<sup>◊</sup> Geraldo F. Oliveira<sup>‡</sup> <sup>°</sup>Univ. of Illinois Urbana-Champaign

Onur Mutlu<sup>‡</sup> <sup>‡</sup>ETH Zürich

We employ the top-down Volcano (Iterator) execution mod





Motivation

Polynesia

. .

Update Propagation . . . . . .

Consistency Mechanism • •

Analytical Engine

1	Introduction	
2	Limitations of HTAP Systems	
3	Polynesia: Overview	
4	Update Propagation Mechanism	
5	Consistency Mechanism	
6	Analytical Engine	
7	Evaluation	
8	Conclusion	
SAFARI Introduction	Motivation Polynesia Update Propagation Consistency Mechanism Analytical Engine Evaluation Conclusion 3	5

# Methodology

- We adapt previous transactional/analytical engines with our new algorithms
  - **DBx1000** for transactional engine
  - C-store for analytical engine
- We use gem5 to simulate Polynesia
  - Available at: https://github.com/CMU-SAFARI/Polynesia
- We compare Polynesia against:
  - Single-Instance-Snapshotting (SI-SI)
  - Single-Instance-MVCC (SI-MVCC)
  - Multiple-Instance + Polynesia's new algorithms (MI+SW) \_
  - **MI+SW+HB: MI+SW** with a 256 GB/s main memory device
  - Ideal-Txn: the peak transactional throughput if transactional workloads run in isolation

SAFARI Introduction . . .

Motivation . . . . . . . .

**Polynesia** • •

Update Propagation . . . . . .

Consistency Mechanism Analytical Engine • •

. . .

Evaluation . . . . . . . .
# End-to-End System Analysis (1/5)



Introduction

SAFAR

Motivation

Polynesia **Update Propagation** . . . . . .

. .

**Consistency Mechanism** . .

**Analytical Engine** 

Evaluation

# End-to-End System Analysis (2/5)



# End-to-End System Analysis (3/5)



# End-to-End System Analysis (4/5)



# End-to-End System Analysis (5/5)



Overall, Polynesia achieves all three properties of HTAP system and has a higher transactional/analytical throughput (1.7x/3.74x) over prior HTAP systems

SAFARI Introduction

Motivation

Polynesia

Update Propagation

n Consistency Mechanism

n Analytical Engine

Evaluation

# **Energy Analysis**

Polynesia consumes 0.4x/0.38x/0.5x the energy of SI-SS/SI-MVCC/MI+SW since Polynesia eliminates a large fraction (30%) of off-chip DRAM accesses



Polynesia is an energy-efficient HTAP system, reducing energy consumption by 48%, on average across prior works

Introduction SAFAR . . .

Motivation . . . . . . . .

**Polynesia** 

• •

Update Propagation . . . . . .

Consistency Mechanism • •

**Analytical Engine** . . .

Evaluation

# More in the Paper

- Real workload analysis
- Effect of the update propagation technique
- Effect of the consistency mechanism
- Effect of the analytical engine
- Effect of the dataset size
- Area Analysis

SAFARI Introduction . . .

Motivation . . . . . . . . **Polynesia** • •

Update Propagation . . . . . .

Consistency Mechanism Analytical Engine • •

. . .

Evaluation . . . . . . . .

# More in the Paper

Real workload analysis

#### Effect of the undate propagation technique

#### Polynesia: Enabling High-Performance and Energy-Efficient Hybrid Transactional/Analytical Databases with Hardware/Software Co-Design

Amirali Boroumand<sup>†</sup> <sup>†</sup>Google

Saugata Ghose<sup>4</sup> Geraldo F. Oliveira<sup>‡</sup> <sup>4</sup>*Univ. of Illinois Urbana-Champaign* <sup>‡</sup>*ETH* 

Effect of the dataset size

Area Analysis



Onur Mutlu<sup>‡</sup>

SAFARI Introduction

Motivation

Polynesia

Update Propagation

Consistency Mechanism

Analytical Engine

Engine Ev

<sup>‡</sup>ETH Zürich

Evaluation

### Outline

1	Introduction
2	Limitations of HTAP Systems
3	Polynesia: Overview
4	<b>Update Propagation Mechanism</b>
5	<b>Consistency Mechanism</b>
6	Analytical Engine
7	Evaluation
8	Conclusion

SAFARI Introduction . . .

Motivation ..... Polynesia • •

Update Propagation .....

• •

Consistency Mechanism Analytical Engine . . .

Evaluation .....

### Conclusion

- Context: Many applications need to perform real-time data analysis using an Hybrid Transactional/Analytical Processing (HTAP) system
  - An ideal HTAP system should have three properties:
    - (1) data freshness and consistency, (2) workload-specific optimization,
    - (3) performance isolation
- **Problem: Prior works cannot achieve all properties of an ideal HTAP system**
- Key Idea: Divide the system into transactional and analytical processing islands
  - **Enables workload-specific optimizations and performance isolation**
- Key Mechanism: Polynesia, a novel hardware/software cooperative design for in-memory HTAP databases
  - Implements custom algorithms and hardware to reduce the costs of data freshness and consistency
  - Exploits PIM for analytical processing to alleviate data movement
- Key Results: Polynesia outperforms three state-of-the-art HTAP systems
  - Average transactional/analytical throughput improvements of 1.7x/3.7x
  - 48% reduction on energy consumption

SAFARI Introduction . . .

Motivation . . . . . . . .

**Polynesia** . .

Update Propagation .....

Consistency Mechanism Analytical Engine • •

. . .

Evaluation

#### **Polynesia:**

Enabling High-Performance and Energy-Efficient Hybrid Transactional/Analytical Databases with Hardware/Software Co-Design

#### Amirali Boroumand Geraldo F. Oliveira

SAFARI Google I

Saugata Ghose Onur Mutlu

ICDE 2022



S ETH Zürich

### Limitations of State-of-the-Art

We extensively study state-of-the-art HTAP systems and observe two key problems:

Data freshness and consistency mechanisms generate a large amount of data movement which causes a drastic reduction in transactional/analytical throughput

They fail to provide performance isolation because of the high resource contention between transactional and analytical workloads

2

#### **Update Gathering & Shipping: Hardware**

To avoid these **bottlenecks**, we design a new hardware accelerator, called update shipping unit



# **Update Application**

Like other relational analytical DBMSs, our analytical engine uses the column-wise data layout and dictionary encoding



#### SAFARI

#### **Update Gathering & Shipping: Algorithm**

Our update shipping algorithm has three major stages:



#### Single-Instance: High Cost of Consistency

Since both analytics and transactions work on the same data concurrently, we need to ensure that the data is consistent

There are two major mechanisms to ensure consistency:

**Snapshotting (Snapshot Isolation)** 



### Snapshotting

#### Several HTAP systems use snapshotting to provide consistency via Snapshot Isolation (SI)



These systems explicitly create snapshots from the most recent version of data and let the analytics run on the snapshot while transactions continue updating data

#### SAFARI

#### Multi-Version Concurrency Control (MVCC)

#### MVCC <u>avoids</u> making full copies of data by keeping several versions of the



When updates happen, MVCC creates a new time-stamped version of data and keeps <u>the old version</u> in a version chain <u>alongside the data</u>

#### **SAFARI**

### **Snapshotting: Drawbacks**

We find that this approach requires frequent snapshot creation to sustain data freshness under high transactional update rate

Two Txn threads Each IM Txn queries Write/read 50%



The overhead comes from Memcpy operation which generates a large amount of data movement and introduces significant interference

#### **More Insights on Data Freshness Challenges**

Our analysis shows that simply providing higher bandwidth (8x) to CPU cores does not address the challenges

We need to take advantage of PIM logic to reduce data movement and resource contention

We find that simply offloading them to general purpose PIM cores does not address the challenges

We need to design custom algorithm and hardware to efficiently execute update shipping/application process

#### Multi-Version Concurrency Control (MVCC)

# We observe that MVCC overhead leads to 42.4% performance loss over zero-cost MVCC



#### We find that long version chains are the root cause of the issue

- Frequent transactional updates create lengthy version chains
- 2 Scan-heavy analytics traverses a lengthy version chain upon accessing a data tuple
- Expensive time-stamp comparison + a very large number of SAFARI random memory accesses

### **Analytical Islands Key Components**

We co-design new algorithms and efficient hardware support for the three key components of an analytical island



#### SAFARI

### **Analytical Islands Key Components**

We co-design new algorithms and efficient hardware support for the three key components of an analytical island



Design an in-memory copy unit that enables highly efficient snapshot creation

### **Analytical Islands Key Components**

We co-design new algorithms and efficient hardware support for the three key components of an analytical island



Simple PIM cores to execute execution engine

# **A Polynesia HW Implementation**

We implement an instance of Polynesia that supports relational transactional and analytical workloads



#### **Consistency Mechanism: Requirements**

Consistency mechanism <u>must not</u> compromise either the throughput of analytical queries or the update propagation rate

**Consistency mechanism has to satisfy two requirements:** 

- Updates must be applied all the time and should not be blocked by analytical queries  $\rightarrow$  Data freshness property
- 2 Analytics must be able to run all the time and should not be blocked by update propagation process → Performance isolation property

# **Analytical Engine: Query Execution**



#### **Analytical Engine: Data Placement**



### **Analytical Engine: Task Scheduler**

For each query, the scheduler makes three key decisions:

- **Decides how many tasks to create**
- 2 Finds how to map these tasks to the available resources (PIM threads)
- **3** Guarantees that dependent tasks are executed in order

### Task Scheduler: Initial Hueristic

Our scheduler heuristic that generates tasks by disassembling the operators of the query plan into operator instances

Query



#### Task Scheduler: Initial Heuristic

We find that this heuristic is **not optimized** for **PIM** and leads to **sub-optimal** performance due to **three reasons**:

The heuristic requires a dedicated runtime component to monitor and assign tasks

- The runtime component must be executed on a general-purpose PIM core
- 2 The heuristic's static mapping is limited to using only the resources available within a single vault group
  - Can lead to performance issues for queries that operate on very large columns
- **3** This heuristic is vulnerable to load imbalance
  - Some PIM threads might finish their tasks sooner and wait idly for straggling threads

# Hueristic

We optimize our heuristic to address these challenges:

We design a pull-based task assignment strategy, where PIM threads cooperatively pull tasks from the task queue at runtime

- We introduce a local task queue for each vault group
- This eliminates the need for a runtime component (first challenge) and allows PIM thread to dynamically load balance (third challenge)

2 We optimize the heuristic to allow for finer-grained tasks

- Partition input tuples into fixed-size segments (i.e., 1000 tuples) and create an operator instance for each partition
- **3** We optimize the heuristic to allow a PIM thread to steal tasks from a remote vault if its local queue is empty
  - This enables us to potentially use all available PIM threads to execute tasks

# Analytical Engine: Hardware Design

Given area and power constraints, it can be difficult to add enough PIM logic to each vault to saturate the available vault bandwidth

Our new data placement strategy and scheduler enables us to expose greater intra-query parallelism

Analytical Island HW Resources



### Wrap up: Single-Instance Systems

While single-instance design enables high data freshness, we find that it suffers from two major challenges:





### **Consistency Mechanism: Algorithm**

Our mechanism relies on a combination of snapshotting and versioning to provide snapshot isolation for analytics

# Our consistency mechanism is based on two key observations:

Updates are applied at a column granularity

Snapshotting a column is cost effective using PIM

#### **Update Propagation: Update Gathering & Shipping**

# Goal: gather updates from transactional threads and ship them to analytical the replica




## Data Freshness Mechanism



### **Data Freshness Mechanism:**

- Update Shipping: gather updates from transactional islands, find the target location in analytical island, and ship them
- Update Application: performs format conversion and applies the update to the analytical replica

Introduction SAFAR . . .

Motivation

Polynesia

• •

Update Propagation . . . . .

**Consistency Mechanism** • •

**Analytical Engine** . . . .

**Evaluation** 

Conclusion

## **Analytical Engine: Hardware**

Given area and power constraints, it can be difficult to add enough

#### **PIM** logic to each vault to saturate the available vault bandwidth

Our new data placement strategy and scheduler enables us to expose greater intra-query parallelism



Introduction . . .

Polynesia • •

**Update Propagation** . . . . .

**Consistency Mechanism** • •

**Analytical Engine** . . . .

**Analytical Engine: Query Execution** Efficient analytical query execution strongly depends on:

Data layout and data placement

## Task scheduling policy

We design a pull-based task assignment strategy, where PIM threads cooperatively pull tasks from the task queue at runtime

### How each physical operator is executed

We employ the top-down Volcano (Iterator) execution model to execute physical operations (e.g., scan, filter, join) while respecting operator's dependencies

SAFARI Introduction . . .

2

Motivation . . . . . . . .

Polynesia • •

Update Propagation ....

Consistency Mechanism • •

Analytical Engine . . . .

Evaluation

Conclusion 75

# End-to-End System Analysis (2/6)



Introduction SAFAR

Motivation

Polynesia

. .

Update Propagation . . . . . .

**Consistency Mechanism** . .

**Analytical Engine** 

Evaluation

Conclusion

6