

# **Open System Config Data:** Firmware to Kernel

Sarathy Jayakumar, Principal Engineer, Intel Corp Natarajan Sivagar, Software Engineer, Google Inc















#### What are we trying to solve

- A rich set of system data that is generated by the BIOS. Examples..
  - **DDR Rank Margining Information**
  - **DDR Post Package Repair information**
  - PCIe Link/Port Margining, Topology etc.
  - eMMC Bus Margining etc
  - Memory Map Topology
  - No standard mechanism to expose this data
- Highly vendor specific. Proprietary mechanisms today



лими













#### How does OS use this data: Examples

- Holistic view aids Data Center operations in efficient decision making
- Examples:
- Failure prediction based on Margining data and error rate
  - Assists in identifying systems that could become unreliable
- FRU Identification to enable faster servicing
  - E.g., Identify which DIMM might need to be replaced
- Margin Drift over time

UMMIT

Health of the FRU over time for planned servicing 









### Mechanisms to Surface Platform Info

- Standard Table based mechanisms
  - SMBIOS, ACPI
- Standard Runtime Mechanisms
  - Platform Runtime Mechanism (Replaces certain categories of SMIs)
  - **UEFI** Runtime Services
- Proprietary Mechanisms (Focus of this Presentation)
  - **BDAT (BIOS Data ACPI Table)**
  - Exposes Platform Info in a Proprietary table and associated Schema













# Challenges with ACPI Tables today

Addresses the lengthy and arduous process of deploying new technologies via ACPI

- Defined by a standards body
- Updates or changes to this requires a lengthy process and needs to be signed off by the WG
- Hence this affects deploying new technologies and features.

#### Seamless update ACPI tables/ data in runtime is not possible.

- OS restrictions and page table mapping and relocation of the ACPI table
- Even if the above is addressed, still require a kernel soft reset for a new table to take effect
  - This is highly disruptive to cloud workloads and affects 99.999% uptime for servers





![](_page_5_Picture_11.jpeg)

![](_page_5_Picture_17.jpeg)

![](_page_5_Picture_20.jpeg)

![](_page_6_Figure_1.jpeg)

![](_page_6_Picture_2.jpeg)

![](_page_6_Picture_5.jpeg)

## Schema brings flexibility and extensibility

- Use richer Redfish / JSON Based schema
  - Describes the published data for interpretation
  - Flexibility Vendor data in an Industry standard format
- Disconnects from ACPI Spec revisions
  - Allows for extensibility
  - New Schemas can be added without Spec updates
- Provides Richer Alternative for ACPI Name Space
  - E.g Can describe System Hierarchy

![](_page_7_Picture_9.jpeg)

лими

![](_page_7_Picture_14.jpeg)

![](_page_7_Picture_16.jpeg)

![](_page_7_Picture_18.jpeg)

# Runtime Update

![](_page_8_Figure_1.jpeg)

![](_page_8_Picture_2.jpeg)

![](_page_8_Figure_3.jpeg)

Publicly Defined Schema (Matched by GUID)

![](_page_8_Picture_5.jpeg)

![](_page_8_Picture_9.jpeg)

![](_page_8_Picture_11.jpeg)

# Example Usage Case -Memory Rank Margining

![](_page_9_Picture_1.jpeg)

![](_page_9_Picture_2.jpeg)

![](_page_9_Picture_3.jpeg)

![](_page_9_Picture_4.jpeg)

![](_page_9_Picture_6.jpeg)

#### How BDAT Fills the gap with DIMM Schema

- Industry standard mechanisms expose platform config to the OS
  - SMBIOS Structures defined by DMTF
  - Various ACPI Tables and Methods
  - JEDEC Specifications on DIMM SPD Data
- The Gap
- BDAT
  - Publish vendor schemas to parse the info provided in the BDAT Table

![](_page_10_Picture_9.jpeg)

![](_page_10_Picture_10.jpeg)

![](_page_10_Picture_12.jpeg)

#### There is no standard mechanism to expose SPD Data, DRAM Timing and Training data

![](_page_10_Picture_16.jpeg)

![](_page_10_Picture_18.jpeg)

## DIMM Serial Peripheral Detect (SPD) Info

![](_page_11_Figure_1.jpeg)

![](_page_11_Picture_10.jpeg)

![](_page_11_Picture_11.jpeg)

• DIMM SPD data in non-volatile storage SPD data is accessed through SMBus • SPD provides relevant info about the DIMM

• Organization, supported timings, serial & product info, manufacturer, etc. • Proprietary maintenance data: Post Package Repair info, fused cell info etc.

Data in SPD is useful during DIMM repair flow

Reading thru SMBus during runtime is not always feasible Should be exposed via standard mechanisms

![](_page_11_Picture_17.jpeg)

![](_page_11_Picture_19.jpeg)

### DIMM Rank Margin Info

#### Rank Margining Tool

![](_page_12_Figure_2.jpeg)

#### Rank Margining Tool provides automated Read/Write DQ and CMD margining

![](_page_12_Figure_4.jpeg)

#### **Per-Bit Margin Interpretation**

	Per	bi	t	П	ar	g	in	s:	F	tx	D
		0	)	-		-	7	8-			-1
	30	8	**	1	8.8	8	Ŕ	19119	9.6	1911	e e
	29	- 94	***	e de	**	t de	ŵ.	90 90	e ste st	n de la	ir ir
	28	10	nin di	1	會會	t the	ġ.	i titi titi	i din di	t din d	i i
	27	Ŕ	**	t th	官官	i St	Ŕ	8.8	1919	1913	e e
	26	10	88	197	R R	r Ar	Ŕ	19119	1919	1913	e e
	25	- 90	89.9	8	**	1	ŵ.	98 98	e ste st	e de s	k sk
	24	8	**	t th	會會	t the	÷.	<b>R</b> A	t the second	t fin i	合合
	23	8	**	19	8.8	R.	Ŕ	8.8	1919	1911	e e
	22	8	88	1	**	÷.	Ŕ	8.8	88	1911	e e
	21	- 90	90 S	r str	97 (A	e de	ŵ.	90 90	e de sé	e de s	ê û
	20	÷.	- 4	t de	會會		÷.	ŵ é	n dê rek	÷÷.	
	19		1	r W	ŵ:		Ŕ	8.8	1919	191	
	18		4	1				ŵ	*	\$	
	17								ŵ.		
	16								會		
	-20										
	-21				ŵ.						
	-22				98 98	9			ŵ.	÷.	
	-23	18	ŵ é	t th	會會	t the		÷.	會	會會	ŝ.
	-24	8	88	19	8.8	÷.		818	8.8	1911	e e
	-25	9	**	r W	R R	9		9 9	**	1911	e e
	-26		19 A	t str	ŵ ŵ	e de	ŵ.	ŵ ŵ	e de sé	n de s	ie de
	-27	- 10	tin di	1	會會	t the	÷.	dir di	i di si	t de s	e e
	-28	8	8.8	r W	8.8	n Mir	Ŕ	8.8	1919	n An I	er ter
	-29	-	88	1	Ŕ Ŕ	*	ŵ.	9 9	**	1913	e e
NO. CO	. D1. A	0:	Ro	<u>(</u> )	Dqs	5	- 1	er?	· b	it	
_0	_1	.2		3		4		5		6	
20	21	18		.8. 12		19 21		20		5	- 1
-43	- 63 -	49	- 4	3		<b>C L</b>			-4	÷.	- 6
36	37	38	3	19	4	40	4	1	4	2	4
21	20	22	- 7	0		20	- 3	1	1	9	2
-25	-24 -	23	-3	7	- 6	23	-3	23	-2	3	-2

- $\bullet$

- $\bullet$

![](_page_12_Picture_11.jpeg)

UMMIT

![](_page_12_Picture_12.jpeg)

![](_page_12_Figure_13.jpeg)

MRC trains the DRAM timings based on SPD data, OEM/Board configuration and margining algorithms MRC programs the memory controller with right timing values MRC disables DIMM/memory channels as required Should be exposed via standard mechanisms

![](_page_12_Picture_15.jpeg)

![](_page_12_Picture_17.jpeg)

### DIMM Software Stack: OS & BMC

BMC thru inband or outband, queries DRAM data for FRU and sensor management.

![](_page_13_Figure_2.jpeg)

![](_page_13_Picture_3.jpeg)

![](_page_13_Picture_4.jpeg)

![](_page_13_Picture_5.jpeg)

SMBIOS & ACPI Tables
AM training data HOB

System Memory Map

- For system memory map the OS relies on standard, Firmware exported data like ACPI tables
- For repair workflow (for DIMM) and other timing info OS/BMC has to use proprietary methods even though underlying technology (DDR3/DDR4 etc) is standard

![](_page_13_Picture_10.jpeg)

**Specifications** 

![](_page_13_Picture_13.jpeg)

![](_page_13_Picture_14.jpeg)

#### Firmware Memory Data Interfaces

![](_page_14_Figure_1.jpeg)

![](_page_14_Picture_2.jpeg)

Proprietary

Firmware

![](_page_14_Picture_4.jpeg)

![](_page_14_Figure_5.jpeg)

![](_page_14_Figure_6.jpeg)

![](_page_14_Picture_7.jpeg)

![](_page_14_Picture_9.jpeg)

# Next Steps

![](_page_15_Picture_1.jpeg)

![](_page_15_Picture_2.jpeg)

![](_page_15_Picture_3.jpeg)

![](_page_15_Picture_4.jpeg)

![](_page_15_Picture_5.jpeg)

![](_page_15_Picture_7.jpeg)

![](_page_15_Picture_8.jpeg)

### ACPI BDAT & Schemas

- Intel's BDAT (BIOS Data ACPI Table) provides a starting point to an open standard
  - Provides Standard ACPI Headers
  - Consists of vendor specific Schemas
- Schemas are basically UUID (GUID) based data structure
- The UUID identifies the type of data associated with the schema
- Multiple schemas will be arranged together to form a coherent data structure
- BDAT schemas provide flexibility in defining new data formats

![](_page_16_Picture_8.jpeg)

![](_page_16_Picture_9.jpeg)

![](_page_16_Picture_10.jpeg)

![](_page_16_Picture_12.jpeg)

![](_page_16_Picture_14.jpeg)

### Proposed Open Firmware Data

- Intel is publishing the BDAT Table and associated schemas in public (link) TBD)
- ACPI Initiatives in place to make it part of the ACPI Specifications
- Schemas are UUID identified and published in the open (not part of ACPI Spec)
- Develop a kernel driver (bdat) to expose BDAT schema
- Create necessary /sys nodes to access BDAT data in userspace

![](_page_17_Picture_6.jpeg)

![](_page_17_Picture_7.jpeg)

![](_page_17_Picture_8.jpeg)

![](_page_17_Picture_9.jpeg)

![](_page_17_Picture_11.jpeg)

#### Call to Action

- Work as a community to make it available in all platforms
- Provide feedback and add missing data
- Create relevant support in Open Source community like standard Linux drivers and tools
- Create additional interfaces like RedFish/IPMI to uniformly identify and utilize this data
- Contacts:
  - OSF Site (<u>https://www.opencompute.org/projects/open-system-firmware</u>)
  - OSF Mailing List (OCP-OSF@OCP-All.groups.io)
  - <u>https://github.com/opencomputeproject/OSF</u>

![](_page_18_Picture_9.jpeg)

![](_page_18_Picture_10.jpeg)

![](_page_18_Picture_14.jpeg)

![](_page_18_Picture_16.jpeg)

![](_page_19_Picture_0.jpeg)

#### SMBIOS Structures

- Memory device (Type 17) structure provide info on individual DIMM such as • DIMM dimensions size, width, speed, voltage etc
  - DIMM identifier such as location string, asset tag, serial number etc
  - Handle to error information structure (Type 18 or 33)
  - Handle to the physical memory array group
  - Usually not filled in in the DIMM is disabled due to training error
- Physical memory array (Type 16) structure groups memory devices to logical unit say a channel or memory riser etc.
- Memory array mapped address (Type 19) & memory device mapped address (Type 20) structures provide ability to translate memory address to DIMM • As SMBIOS allows variable structure size, vendors use the extra data for storing
- proprietary info about the DIMM

![](_page_20_Picture_9.jpeg)

![](_page_20_Picture_10.jpeg)

![](_page_20_Picture_12.jpeg)

![](_page_20_Picture_13.jpeg)

![](_page_20_Picture_15.jpeg)

![](_page_20_Picture_16.jpeg)

#### ACPI Tables

- Static Resource Affinity Table (SRAT)
  - Describes the resource affinity domains in the system
  - (NUMA) of the system
- System Locality Distance Information Table (SLIT) • Describes the locality of NUMA (Non-Unified Memory Architecture)
  - nodes
  - The distance is given in terms of Memory Latency

![](_page_21_Picture_7.jpeg)

![](_page_21_Picture_8.jpeg)

Critical structure to understand the Non-Unified Memory Architecture

![](_page_21_Picture_10.jpeg)

![](_page_21_Picture_14.jpeg)