

## **Verification Futures**

July 2025

SRAM: Dead or alive?



#### What is an SRAM?



- + Volatile



# Six transistors per bit cell





### SRAM: History and timeline



#### **chiplet**<sup>®</sup>

## Why is it so ubiquitous?

- 1. Best read/write latencies
- 2. Fully deterministic
- 3. Ease of integration with logic
  - 1. SRAM memory compiler for every process node
  - 2. Generate any shape and the size you want for your design

<u>Memory</u>	Read latency	<b>Density</b>	
SRAM	0.5-2ns	~60-120F2	Available on every process node as part of logic foundries SRAM compilers.
DRAM	10-15ns	~6-8F2	Special foundries that manufacture DRAM chips - currently cannot be integrated together with logic which adds additional off-chip latencies in the order of 50-100ns.
MRAM	10-35ns	~20-40F2	Very limited foundries and process nodes.
ReRAM	10-100ns	~4F2	Limited availability: 22nm FDSOI, 130nm etc

*F*<sup>2</sup> represents feature size relative to process node



### Example designs

## Die shots of microprocessors using a large quantity of SRAMs. Some of the smaller macros not shown.



Sandy Bridge [32nm]



Power6 [65nm]



McKinley [180nm]

Source: https://en.wikichip.org/wiki/static\_random-access\_memorygoogle\_vignette

**chiplet**<sup>1</sup>

#### However!



Source: <u>H. Chang et al.</u>

SRAM scaling has *significantly slowed* down below 5nm in density and performance

# SRAM uses large chip area which could be used for logic



Source: AMD Ryzen 5000 Zen 3 Desktop CPU Gets First High-Res Infrared Die Shot, Vermeer Fully Detailed

#### chiplet<sup>®</sup>

#### Solution?

- 1. When large arrays of SRAMs are needed, put them on a separate die
- 2. Vertically stack them on top of compute
  - 1. Why? Connections via silicon interposers require SerDes latencies in the range of 50-100ns





### Why now?

- 1. Designs are hitting reticle limits
- 2. Advanced 3D integration with very high yields
  - 1. Bumpless sub 10um pitch IOs in mass production for several years
- 3. Wafer costs disparity
- 4. Reliability/scaling issues

Let's dive in a bit more.



#### **Break-even**

3 break-even factors to consider for switching to 3D stacked SRAMs.





#### **Break-even: PHY area**

PHY area > SRAM macro

# PHY area depends on the number of IOs, process node, the type of bonding and the pitch between the IO cells.

#### 1MB

Hybrid bond, 9um pitch, 300 IO cells at 7nm





#### **Break-even: Costs**



Volume production costs for a 128MB SRAM chiplet

Doesn't factor the exorbitant NRE costs!

Break-even for costs is very subjective.



#### **Break-even: Performance**

Need to measure system level performance. Factors to consider include Fmax improvements due to smaller compute die, lower latencies within the compute die, higher latencies to off chip SRAMs and higher SRAM capacity.

chiplet<sup>®</sup>

### Example: Al inference

Al Model	Men	nory	#Chips	<b>Optical cables</b>	
size	Size	Stack		#Number	Miles
7B	16GB	2D	70	315	1.9
		3D	16	70	0.4
	20TB	2D	95,325	428962	2659
10T		3D	16,850	75825	470

A memory stacked approach can save 4-10x number of chip modules or optical cables for the same AI model.





#### SRAM chiplets: configurations

Products	Small	Medium	Large
Memory	16MB	64MB	128MB
Chiplet Area	4.5mm <sup>2</sup>	16.8mm <sup>2</sup>	33.7mm <sup>2</sup>
I/Os	256	512	1024
Macro area	0.02mm <sup>2</sup>	0.05mm <sup>2</sup>	0.05mm <sup>2</sup>
Bandwidth	128Gb/s	256Gb/s	512Gb/s
Process		7nm	
Bond pitch		9um	





## Thank you

kauser.johar@chipletti.com