



AI-ASIC Engagement Models: Alchip vs GUC in the Race Beyond GPUs

Author: Lai Yit Loong

Date: August 31, 2025

Executive Summary

Graphics Processing Units (GPUs) set off the AI revolution, but the very success of AI has exposed their limits. Rising cost per token processed, excessive power consumed per rack, supply bottlenecks, and mismatch between workloads all highlight the need for alternatives.

That alternative is now arriving in the form of **AI-ASICs (Application-Specific Integrated Circuits)** and other custom accelerators. These purpose-built chips trade generality for efficiency, achieving better performance per dollar, better performance per watt, and faster time-to-capacity.

This report explores how custom silicon is reshaping compute economics, focusing on two key design-service providers: **Alchip** and **GUC (Global Unichip Corp.)**. Their different engagement models—whether providing low-margin logistics services or high-value design platforms—determine both their earnings power and their long-term strategic position.

Key insights include:

- **Why alternatives matter:** Training and serving AI models are different tasks. Custom silicon provides efficiency by tailoring hardware to the specific needs of each.
- **Alchip:** Building a reputation as an “open-arms prime contractor,” offering complete 2 nanometer (nm) design platforms, chiplet integration, and rack-level connectivity. Anchored by a multi-billion-dollar 3nm hyperscaler program and diversifying into automotive chips.

- **GUC:** A close partner of TSMC, scaling turnkey projects for cloud customers but currently trapped in thin-margin logistics work (TK3). The challenge is to climb to higher engagement levels (TK2/TK1) where more value is captured.
- **Industry fulcrum:** The rule is clear — *Engagement model equals margin model*. Revenue growth alone does not guarantee profit leverage unless the service provider owns more of the design, packaging, and yield process.

Takeaway: In the AI-ASIC race, allocation of scarce resources is strategy, packaging is becoming the product, and turnkey without leverage risks becoming a commodity.

Table of Contents

1. Introduction
 2. Why the World Needs More Than GPUs
 3. Market Dynamics
 - 3.1 Hyperscaler Programs (3nm → 2nm)
 - 3.2 Automotive ADAS (Advanced Driver Assistance Systems)
 - 3.3 Crypto Demand vs Capacity Allocation
 4. Engagement Models: TK3 vs TK2 vs TK1
 5. Technology Trends: 3.5D as the Middle Ground
 6. Supply Chain and Standards
 7. Company Deep Dives
 - 7.1 Alchip
 - 7.2 GUC
 8. Key Numbers
 9. Strategic Analysis
 10. Impact and Recommendations
 11. Glossary
-

1. Introduction

The transition from general-purpose compute to specialized accelerators is the defining shift of this AI cycle. GPUs powered the early wave, but as models scale larger and serving workloads dominate cost structures, their generality turns into inefficiency.

Custom silicon—AI-ASICs, low-latency processors like Groq's LPU, and wafer-scale engines like Cerebras's WSE—are emerging as the economic and strategic response.

Taiwan's **Alchip** and **GUC** represent two archetypes:

- **Alchip** focuses on building alliances, platforms, and differentiated design capability.
- **GUC** leverages its close alignment with TSMC, the world's most advanced foundry, to scale turnkey projects.

Both are racing to capture hyperscaler budgets, expand into automotive markets, and secure scarce capacity in high-bandwidth memory (HBM) and advanced packaging.

2. Why the World Needs More Than GPUs

GPUs ignited the AI boom, but scale exposed limits:

- **Economic limits:** The cost per token processed by large language models continues to rise. GPU racks also consume large amounts of electricity, raising both operating costs and sustainability concerns.
- **Physical limits:** Larger models push against bandwidth ceilings, latency bottlenecks, and thermal issues. Monolithic GPUs cannot scale efficiently.
- **Supply limits:** GPUs rely on leading-edge manufacturing capacity, which is rationed by foundries. Long lead times mean hyperscalers wait in line for allocation.
- **Mismatch of tasks:** Training a large model is not the same as serving billions of user queries. Training favors flexibility; serving demands efficiency and predictability.

Why AI-ASICs now:

- **AWS Trainium2** delivers both training and inference performance with high-bandwidth memory (96 GiB HBM, 2.9 TB/s bandwidth).
- **Google TPU Trillium** delivers four times the training performance and up to three times inference throughput compared to the previous generation.
- **Microsoft Maia 100** integrates silicon, networking, and cooling for Azure workloads.
- **Specialists like Groq (LPU)** target ultra-low latency, while **Cerebras (WSE-3)** exploits wafer-scale silicon for massive on-chip memory.

Supply chain threads and who plays where:

1. **Foundry and advanced packaging:** TSMC, Samsung, Intel provide leading-edge wafers and CoWoS/SolC packaging.
 - **Alchip and GUC:** dependent on TSMC allocation; GUC enjoys equity alignment.
 - **Broadcom, MediaTek, Marvell:** fabless giants with scale leverage for cloud ASICs.
2. **High-bandwidth memory (HBM):** SK hynix, Samsung, Micron are key suppliers.
 - Memory is the throttle; without it, compute capacity stalls.
3. **Substrates and OSAT (Outsourced Assembly and Test):** Companies like Ibiden, Unimicron, ASE, and Amkor are critical.
 - **Alchip:** pre-books OSAT capacity through partnerships.
 - **GUC:** executes turnkey assembly at scale, but margins remain thin.
 - **Broadcom/MediaTek:** secure priority through volume demand.
4. **Die-to-die standards (UCle 3.0):** Provides interoperability between chiplets at speeds up to 64 GT/s.
 - **Alchip:** designs chiplet-ready platforms at 2nm.
 - **Broadcom/Marvell:** push standards for network-compute integration.

- **MediaTek:** focuses on consumer and edge AI ecosystems.

Thesis: Diversity of compute is not optional — it is resilience. Heterogeneous accelerators, packaged coherently, will be required to keep up with model growth and user demand.

3. Market Dynamics

3.1 Hyperscaler Programs

- **Alchip:** Anchored by a \$4–6 billion 3nm accelerator program (shipments 2026–2028). Already preparing a 2nm successor platform that mixes advanced compute with cost-optimized I/O chiplets. Competing in a very complex 3.5D AI project against tier-one competitors like Broadcom, MediaTek, and Marvell.
- **GUC:** Supporting CPU ramp in 4Q25 and AI-ASIC ramp in 2026. Revenue could grow more than 40% year-over-year, but much of this sits in TK3 (logistics-heavy) engagements, which cap margins at 3–5%.

3.2 Automotive ADAS (Advanced Driver Assistance Systems)

- **Alchip:** Won a design program with Li Auto, expected to contribute \$300–400 million in 2026 with gross margins above 25%. Represents diversification beyond hyperscalers into automotive.

3.3 Crypto

- **Alchip:** Contributes \$70–100 million annually, but upside is limited because capacity is prioritized for AI programs.
 - **GUC:** Sees strong crypto demand, but supply is bottlenecked at foundries and OSATs.
-

4. Engagement Models: TK3 vs TK2 vs TK1

- **TK3 (Turnkey Level 3):** Primarily logistics and yield tuning. Vendor ensures chips flow through the foundry and assembly lines, but contributes limited

unique IP. Margins are thin, typically 3–5%. This is where GUC is heavily concentrated today.

- **TK2 (Turnkey Level 2):** Involves full physical design responsibilities (place-and-route, power, timing closure). Adds more value and allows higher gross margins.
- **TK1 (Turnkey Level 1):** End-to-end responsibility, including front-end design and architecture (RTL to GDSII). Highest margin and highest risk. Alchip is climbing toward TK1 through NRE-intensive hyperscaler projects.

Key insight: The engagement model directly dictates the margin model. Without moving beyond TK3, revenue growth does not translate into earnings leverage.

5. Technology Trends: 3.5D as the Middle Ground

- Monolithic scaling of chips beyond ~800 mm² is impractical due to yield collapse.
 - **3.5D packaging** blends chiplets, stacked dies, and interposers to strike a balance between performance, cost, and manufacturability.
 - Allows node mixing: 2nm compute logic paired with 3nm/5nm I/O and memory controllers.
 - In this era, **packaging is no longer just assembly — it is the product.** Yield learning becomes a competitive moat.
-

6. Supply Chain and Standards

- **TSMC:** The gravitational center. Allocation of N3/N2 wafers and CoWoS/SolC packaging defines winners.
- **HBM (High-Bandwidth Memory):** The true bottleneck. Without it, accelerators cannot be fully utilized.
- **OSATs (Outsourced Assembly and Test):** ASE, Amkor, and others now control throughput bottlenecks.
- **ABF substrates:** Scarcity constrains delivery across the ecosystem.

- **UCle 3.0:** Standard chiplet bus enabling multi-vendor interoperability and reducing lock-in.
-

7. Company Deep Dives

7.1 Alchip – The Open-Arms Prime

- Anchored by multi-billion-dollar 3nm hyperscaler program; successor 2nm platform already underway.
- Alliance with Astera Labs allows inclusion of rack-scale connectivity solutions (PCIe, CXL retimers, memory expansion).
- Diversifying into automotive with Li Auto ADAS project.
- **Strengths:** Modularity, speed of proposal, alliances.
- **Risks:** Scarce HBM allocation, capital intensity, and export restrictions on China programs.

7.2 GUC – The TSMC-Aligned Integrator

- TSMC is its largest shareholder and sole foundry partner. This provides privileged access but also dependence.
 - Scaling turnkey hyperscaler programs at volume.
 - Revenue robust, but concentrated in TK3 logistics-heavy work with thin margins.
 - **Unlock:** Climbing to TK2/TK1 by owning package bring-up, test, and yield analytics.
 - **Risks:** Heavy customer concentration in North America, volatility in crypto, and limited earnings leverage.
-

8. Key Numbers: Financial Comparison – Alchip vs GUC (USD)

Metric	Alchip Technologies	Global Unichip Corp. (GUC)
Market Cap (USD B)	~10.0	~5.5
Revenue (TTM, USD B)	~1.47	~0.80
Net Income (TTM, USD B)	~0.20	~0.11
Gross Margin	~21% – diversified (hyperscaler ASICs + higher-margin auto/ADAS)	~32.5% – strong cost efficiency but turnkey-heavy mix
Operating Margin	~16% – supported by auto diversification and NRE-heavy projects	~16.5% – steady but capped by logistics scope
Net Profit Margin	~13.5% – improved by diversification into auto and custom design scope	~14.1% – relatively stable but limited by turnkey model
Free Cash Flow (USD B)	~0.48 – strong cash generation, net cash positive, minimal debt	Not disclosed publicly
2026 Revenue Growth Outlook	Revenue expected to boost from hyperscaler program + auto ADAS	~+40% YoY, driven by CPU + AI-ASIC ramps, still TK3-heavy
2026 EPS Growth Outlook	Sharp ramp from ADAS margins and 2nm ASIC launches	+8–9%, capped by low-margin turnkey exposure

9. Strategic Analysis

Dimension	Alchip	GUC
Revenue Scale	Multi-billion 3nm/2nm programs	+40% YoY 2026 growth, but TK3-heavy
Margin Profile	25%+ auto programs, NRE leverage	3–5% gross margin on turnkey
Alliances	Astera Labs, IP vendors, modular system	Tight alignment with TSMC
Technology Edge	2nm platform + 3.5D readiness	Operational scale in turnkey

Dimension	Alchip	GUC
Key Risks	Allocation, geopolitics, capital needs	Mix skew, hyperscaler concentration

10. Impact and Recommendations

For CEOs of design houses and fabless firms:

1. Separate roadmaps for training and serving — one size no longer fits all.
2. Pre-book scarce resources early: high-bandwidth memory, CoWoS/SolC packaging, ABF substrates.
3. Standardize interfaces (UCIe) to preserve flexibility in vendor selection.
4. Productize “design-ops”: make yield, test, and reliability part of the service offering and price for delivered outcomes.
5. Build alliance gravity — connectivity, firmware, and memory partnerships — so you sell complete rack-level solutions, not just chips.

For Investors:

- Favor vendors that secure allocation rights, build alliances, and move up the engagement curve (TK3 → TK2/TK1).
- Look beyond revenue growth — focus on yield-to-ship metrics as the true driver of profitability.

For Policymakers:

- Recognize that OSATs and packaging are as strategic as fabs.
- Extend incentives and national security frameworks to substrates, interposers, and HBM fabs.

11. Glossary

- **AI-ASIC:** Application-Specific Integrated Circuit for AI workloads.
- **TK3/TK2/TK1:** Levels of turnkey engagement. TK3 = logistics support, TK2 = full physical design, TK1 = full architecture and RTL.

- **NRE (Non-Recurring Engineering):** One-time design fees paid upfront for complex projects.
- **SoIC (System on Integrated Chips):** TSMC's advanced 3D stacking technology.
- **CPO (Co-Packaged Optics):** Integration of optical I/O inside chip packages.
- **UCIe (Universal Chiplet Interconnect Express):** Industry standard for die-to-die communication between chiplets.
- **HBM (High-Bandwidth Memory):** Stacked memory technology with very high throughput.
- **ADAS (Advanced Driver Assistance Systems):** Automotive electronics enabling driver safety features.
- **OSAT (Outsourced Semiconductor Assembly and Test):** Companies providing packaging, assembly, and final testing of chips.