*Original Article*

# Small Language Models and Neuro-Symbolic AI in Zonal Architectures: The Rise of Small Language Models (SLMs) in Constrained Environments

Naresh Kalimuthu
Indepentent Researcher, USA.

*Abstract - The automotive industry faces a pivotal moment with the rise of Software-Defined Vehicles (SDVs), Zonal Electronic/Electrical (E/E) Architectures, and Generative AI. As vehicle architectures shift from centralized domain structures to decentralized zonal networks to handle more complex wiring and data flow, the need for local intelligence grows. However, deploying Large Language Models (LLMs) on resource-limited vehicle edge devices is challenging due to high computational latency, energy demands, and safety concerns. This paper explores the potential of Small Language Models (SLMs)—under 7 billion parameters, such as Microsoft's Phi-3 and TinyLlama—as a way to embed advanced reasoning directly into Zonal Controllers. It also tackles the core issue of balancing the unpredictability of generative AI with the strict determinism required by safety standards such as ISO 26262. We suggest and evaluate Neuro-Symbolic AI as an essential layer of the architecture, using symbolic logic to verify neural network outputs in real time. By examining hardware options (NXP S32, TI TDA4), virtualization tools (Adaptive AUTOSAR), and safety protocols (RoboGuard, SYNAPSE), this study shows that SLMs, guided by symbolic logic, can support resilient, decentralized vehicle autonomy, reduce dependence on cloud connectivity, and uphold safety integrity.*

*Keywords - Small Language Models (Slms), Zonal Architecture, Neuro-Symbolic AI, Software Defined Vehicle (SDV), Edge Computing, Adaptive AUTOSAR, Deterministic Safety, Real-Time Systems, Knowledge Distillation, ISO 26262, Time-Sensitive Networking (TSN).*

## 1. Introduction

The modern car is quickly transitioning from a simple mechanical device to a highly sophisticated, interconnected platform at the edge of computing technology. This shift, often called the "Software-Defined Vehicle" (SDV), is driven by consumer demand for richer digital experiences, advanced driver-assistance systems (ADAS), and autonomous driving features. However, this surge in functionality has overwhelmed traditional automotive Electronic/Electrical (E/E) architectures. Old-style distributed systems, which link each function directly to its own Electronic Control Unit (ECU), have led to vehicles with over 100 ECUs and wiring harnesses weighing as much as 60 kilograms, posing serious challenges in manufacturing, weight management, and software maintenance.

To address these physical and logical bottlenecks, the industry is moving quickly toward Zonal Architectures. In this approach, the vehicle is divided by physical location such as the front-left or rear-right zones rather than by function, such as powertrain or infotainment. High-performance Central Compute Units (CCUs) handle complex processing tasks, while "Zonal Controllers" gather sensor data and manage power distribution within their designated zones. This transition is expected to cut wiring complexity by up to 30% and make it easier to add new features through Over-the-Air (OTA) updates.

Alongside this hardware revolution, Generative Artificial Intelligence is also rising. Large Language Models (LLMs) have shown remarkable abilities in reasoning, understanding natural language, and generating code. However, deploying "foundation models" with over 100 billion parameters is generally incompatible with automotive edge constraints. Vehicles have limited power budgets, need effective thermal management, and depend on cloud connectivity for inference, which would cause unacceptable latency and privacy concerns.

This contrast the requirement for sophisticated AI reasoning versus the limitations of edge devices—has driven the emergence of Small Language Models (SLMs). Examples such as Microsoft's Phi-3, Google's Gemma, and the open-source TinyLlama demonstrate that effective reasoning is possible with fewer than 4 billion parameters, suggesting that they can be deployed on automotive-grade silicon in principle.

However, a significant obstacle is safety: automotive systems must adhere to stringent functional safety standards such as ISO 26262, which require determinism, traceability, and "freedom from interference." Neural networks are inherently probabilistic, making them susceptible to "hallucinations"—plausible yet incorrect outputs—which are unacceptable in safety-critical situations. A braking system, for example, cannot rely on a "likely" correct decision; it must function based on absolute certainty.

This paper suggests that the answer resides in Neuro-Symbolic AI, a combined approach that integrates the learning and perceptual strengths of neural networks with the precise, rule-based reasoning of symbolic logic. By integrating SLMs into Zonal Controllers and incorporating symbolic guardrails, OEMs can achieve localized, intelligent, and safe vehicle autonomy. This paper discusses the technical hurdles, architectural strategies, and real-world testing for this integration.

## 2. Research Topics

Integrating probabilistic SLMs into safety-critical, resource-limited zonal architectures involves complex engineering challenges. This research addresses three key topics that constitute the main technical hurdles to this technology's adoption.

### 2.1. Hardware-Software Co-Design and Resource Contention in Mixed-Criticality Zonal Controllers

The Architectural Conflict: In a pure Zonal Architecture, the Zonal Controller (ZC) serves as a dedicated gateway and power-distribution hub. As the SDV concept evolves, the ZC assumes more responsibility for "Edge Processing" to relieve the central computing system. This creates a mixed-criticality workload environment. On one side, the ZC must perform strict real-time control tasks such as managing smart fuses (eFuses) to disconnect short circuits within microseconds or activating local actuators for braking and steering. These operations usually run on Real-Time Operating Systems (RTOS) like AUTOSAR Classic, which demand deterministic execution with sub-millisecond precision jitter.

On the other hand, deploying an SLM (such as a 2-billion-parameter model for local voice diagnostics or sensor anomaly detection) creates a bursty, compute-intensive workload. Even quantized SLMs require substantial memory bandwidth (RAM) and rely on large vector processing units.

### 2.1.1. The Research Challenge: Resource Contention

The main research challenge is Resource Contention. How can one System-on-Chip (SoC) run a memory-intensive, non-deterministic SLM alongside a safety-critical control loop without violating the "Freedom from Interference" (FFI) requirement of ISO 26262 ASIL-D?

If the SLM saturates the internal interconnect fabric or empties the shared Last-Level Cache (LLC), it may cause stalls in the real-time cores, resulting in missed safety task deadlines. This issue, called "noisy neighbor" interference, poses a significant challenge to safety certification.

- Hardware Heterogeneity Analysis: Present automotive silicon solutions, such as the NXP S32G and Texas Instruments TDA4VM, address this by using heterogeneous multi-core architectures.
- Application Cores: High-performance Arm Cortex-A53 or A72 cores are allocated to the "Performance Domain," running POSIX-compliant OSs such as Linux or QNX to host the SLM.

- Safety Cores: Lockstep Arm Cortex-R52 or Cortex-M7 cores serve as the "Safety Domain" and execute the RTOS.

However, physically isolating cores does not ensure that shared resources like DRAM controllers and system buses are also separated. Research should explore how effective hardware-enforced isolation mechanisms such as Quality of Service (QoS) regulators on the interconnect and System Memory Protection Units (SMPUs) are specifically for SLM inference workloads.

### 2.2. The Stochastic-Deterministic Conflict: Ensuring Safety via Neuro-Symbolic AI

The Safety Paradox: The automotive industry is based on the principle of determinism, where a safety system must consistently produce the same output (actuation) given the same inputs (sensor readings). Both functional safety standards (ISO 26262) and Safety of the Intended Functionality (SOTIF / ISO 21448) depend on traceability and validation.

Generative AI, including SLMs, functions based on probabilistic models. An SLM predicts the next token in a sequence using a statistical distribution derived from extensive datasets. This approach inherently introduces variability and, importantly, the risk of hallucinationthe production of factually incorrect or physically impossible outputs. For instance, an SLM serving as a "digital co-pilot" could misinterpret a sensor code and mistakenly advise the driver to continue driving despite a severe battery thermal runaway event.

### 2.2.1. The Research Challenge: bridge the gap between Neural Probabilities and Symbolic Certainties

The challenge is to connect Neural Probabilities with Symbolic Certainties. We can't just "trust" the neural network; instead, we need a mechanism that can formally verify the SLM's output in real-time before acting on it.

This points to Neuro-Symbolic AI. Pure symbolic AI (rule-based systems) is fully deterministic and verifiable but fragile; it cannot handle the unstructured nuances of natural language or complex visual scenes. Pure Neural AI (Deep Learning) is robust and adaptable but opaque and unreliable. The research focuses on defining the architecture of a Hybrid System specifically for the automotive context where the SLM manages perception and interaction, while a Symbolic Engine enforces strict logical constraints (e.g., temporal logic specifications).

We need to explore how to convert automotive safety manuals and regulations, such as Euro NCAP rules, into machine-readable symbolic logic like Linear Temporal Logic or Answer Set Programming. These can then function as a real-time "supervisor" during operation SLM.

## 2.3. Bandwidth Saturation and the Optimization of Distributed Inference

- The Data Deluge: One of the main aims of Zonal Architecture is to reduce weight and complexity, but another crucial reason is efficient data handling. Modern autonomous vehicles generate between 4 and 20 terabytes of data every day. In a centralized system, all this raw data such as LiDAR point clouds and 4K videos needs to be transmitted to the central computer for processing.
- The Bandwidth Bottleneck: Even with the implementation of Automotive Ethernet (1000BASE-T1) and Time-Sensitive Networking (TSN), the backbone bandwidth remains limited. Overloading the backbone with raw sensor data can increase latency and jitter, which may hinder timely critical control response messages.

### 2.3.1. The Research Challenge: Distributed Inference Optimization

This presents a challenge in Distributed Inference Optimization. Instead of sending raw data, can the Zonal Controller leverage an embedded SLM to execute "Semantic Compression"?

Scenario: A rear-zone camera records an image of a sign.
Raw Data: 500 MB/s video stream.
Semantic Data: The text 'Speed Limit 50'.

The challenge is identifying the optimal "Split Point." Which SLMs are capable of efficiently running on the limited hardware of a Zonal Controller (typically 2-10 TOPS) versus the Central Compute (200+ TOPS)? We need to benchmark models like TinyLlama (1.1B) and Phi-3 Mini (3.8B) against the hardware accelerators available in zonal SoCs to assess the viability of this distributed intelligence. Additionally, maintaining consistency between distributed SLMs where different zones, such as the front-left and front-right, might infer different contexts adds complexity to the vehicle's overall "World Model."

## 3. Recommendations / Mitigation Strategies

To effectively incorporate SLMs into zonal architectures within safety and resource limits, a multi-layered approach that includes Neuro-Symbolic guardrails, comprehensive model optimization, and strict architectural isolation is required.

### 3.1. Implementation of Type-2 Neuro-Symbolic Guardrails ("The Sandwich Model")

To mitigate the non-determinism of SLMs, OEMs must adopt a Type-2 Neuro-Symbolic architecture, often referred to as the "Sandwich Model" in safety-critical AI research. This architecture wraps the probabilistic neural model between two layers of deterministic symbolic processing.

- Layer 1: Symbolic Pre-Processing (Grounding): Before the user query or sensor data reaches the SLM, it is "grounded" via a symbolic processor. This process incorporates context from the vehicle's verifiable state such as speed, gear status, and error codes into the prompt. Doing so narrows the search space for the SLM and grounds its reasoning in the current state of information.
- Layer 2: Neural Inference (The SLM):The SLM (such as a fine-tuned Phi-3) analyzes the grounded input to produce a response or action plan, utilizing the model's semantic comprehension flexibility.
- Layer 3: Symbolic Post-Processing (Verification):

The SLM output is not directed to the driver or actuator. Instead, it is forwarded to a Symbolic Logic Engine, which incorporates a "World Model" encoded with strict constraints based on the ISO 26262 safety standard.

**Example Rule:** IF Vehicle_Speed > 0 THEN Door_Lock_Command != UNLOCK.

- Verification: The engine checks the SLM's output. If the SLM incorrectly suggests unlocking the door while moving (a hallucination), the Symbolic Engine blocks the command and activates a deterministic fallback routine, such as "Command rejected: Safety Violation.'
- Recommendation: Using formal verification techniques such as Linear Temporal Logic (LTL) enables a mathematical proof that the system always complies with safety constraints, regardless of the AI's behavior.

### 3.2. Knowledge Distillation and Quantization for "Edge-Native" Models

Running a standard 7B or 70B parameter model on a Zonal Controller isn't feasible. To implement intelligence at the "Edge," an aggressive compression strategy is necessary.

- Knowledge Distillation (KD): Use a powerful 'Teacher' model like GPT-4 or Llama-3-70B to train a smaller 'Student' SLM, such as TinyLlama-1.1B. The Teacher creates synthetic datasets covering automotive diagnostics, rare driving scenarios, and technical manual interpretations. The Student is trained specifically on this data.
- Benefit: The student ignores irrelevant details, like poem writing, and overfits to the automotive domain. It maintains high accuracy for the specific use case while lowering the parameter count 10-50x.
- Quantization: Transition from Floating Point (FP16/FP32) to Integer math (INT4 or INT2). Research shows that modern SLMs can preserve over 90% of their reasoning ability at 4-bit precision quantization.

Impact: A 3B parameter model in FP16 format needs about 6GB of RAM, while in INT4 format, it requires roughly 1.5GB. This makes it compatible with the memory capacity of automotive Zonal Processors such as the NXP S32G.

### 3.3. Hypervisor-Based Isolation and Traffic Shaping

To address the resource contention problem on the Zonal Controller SoC, it is essential to implement strict hardware virtualization.

- *Hypervisor:* Use a bare-metal hypervisor like QNX Hypervisor or Xen to divide the hardware resources.
- Partition A (Safety): This partition, dedicated to the AUTOSAR Classic stack for strict real-time control, is allocated the highest priority for CPU cycles and interconnect access bandwidth.
- Partition B (Intelligence): Runs a guest OS (Linux) that includes the SLM container.
- Memory Protection: Configure the System Memory Protection Unit (SMPU) to physically prevent Partition B from accessing the memory address ranges of Partition A. This helps avoid "memory leaks" or "buffer overflows" in the AI stack from causing safety issues.
- Network Isolation (TSN): Implement IEEE 802.1Qbv (Time-Aware Shaper) on the Ethernet interface, which divides the network schedule into designated time slots windows.
- Window 1: Designated specifically for Safety Critical Control traffic, such as Brake signals.
- Window 2: Available for AI/SLM data traffic, this guarantees that even if the SLM produces a data burst, it cannot delay the transmission of a brake command by even a microsecond.

## 4. Recommendations and Goals Achieved Based on Case Studies

The integration of SLMs and Zonal Architectures is transitioning from theoretical concepts to practical application. The case studies that follow demonstrate how these mitigation strategies have been implemented to produce measurable enhancements in efficiency, safety, and overall performance capability.

### 4.1. Case Study 1: The "RoboGuard" Architecture – Validating Neuro-Symbolic Safety

- Context: Research led by Ravichandran et al. (2025) introduced RoboGuard, a two-stage guardrail system designed to protect LLM-enabled robots from unsafe behaviors and adversarial attacks. Although initially created for general robotics, the architecture is readily applicable to autonomous vehicles actuators.
- Implementation: The system used a "Root-of-Trust" LLM to translate environmental context into formal constraints, with a control synthesizer acting as the symbolic gatekeeper.
- Goal Achieved: During simulations involving "jailbreaking" attacks where the AI was deceived into unsafe behaviors RoboGuard reduced the execution of unsafe plans from 92% to below 2.5%.
- Automotive Improvement: This case study confirms the effectiveness of the Neuro-Symbolic mitigation strategy. It demonstrates that a "Logical Guardrail" is more than just a theoretical idea; it is a practical

software component that filters a generative model's output in real-time, preventing physical harm even if the AI errs. This offers a pathway for certifying AI systems under ISO standards 26262/SOTIF.

### 4.2. Case Study 2: JAL & Microsoft Phi – The Viability of "Offline" Intelligence

- Context: Japan Airlines (JAL) implemented Microsoft's Phi-4 model to allow cabin attendants to create detailed reports and queries directly on mobile tablets without needing cloud connectivity.
- Implementation: The SLM was optimized for local execution on the device (Edge AI), enabling it to perform natural-language summarization and information retrieval tasks that previously depended on an server-side LLM.
- Goal Achieved: The deployment demonstrated that a 3.8B-parameter model could achieve reasoning capabilities comparable to those of much larger models (such as Mixtral 8x7B) on specific tasks. It reduced report generation time while maintaining 100% data privacy (no data left the device).
- Automotive Improvement: This supports the Distributed Inference approach directly. It demonstrates that a Zonal Controller, with comparable computing power to a high-end tablet, can run a capable SLM. This makes features such as "Offline Voice Assistant" or "Local Diagnostics' possible in vehicles, providing functionality even in tunnels or dead zones, and lowering data transmission costs by processing text locally.

### 4.3. Case Study 3: The SYNAPSE Framework – Determinism in Engineering

- Context: The SYNAPSE (Symbolic Neural Architecture for Predictive Structural Engineering) project developed a chatbot for structural analysis a domain where "hallucinated" calculations can lead to building collapse.
- Implementation:The system employed a hybrid approach, with an AI model managing the user's natural language queries, while all essential engineering calculations were carried out by a deterministic symbolic engine (3Muri).
- Goal Achieved: The hybrid system reached 94% accuracy with response times under 2 seconds. Importantly, the symbolic layer guaranteed that none of the outputs ever breached Eurocode building regulations, regardless of the AI's involvement "creativity".
- Automotive Improvement: This confirms that Neuro-Symbolic AI can be effectively used for Vehicle Diagnostics. An SLM could interpret a driver's vague issue, like "The car feels shaky," while the diagnostic logic would remain consistent with the vehicle's precise telemetry data, ensuring the AI doesn't recommend repairs that contradict the actual sensor readings.

### 4.4. Case Study 4: TinyLlama Benchmarking on Edge Silicon

- Context: Benchmarking the TinyLlama (1.1B) model on resource-limited embedded hardware to evaluate its suitability for "tiny" edge applications devices.
- Implementation: Researchers evaluated the model's throughput and memory usage across different hardware configurations, focusing on optimizing for 4-bit and 8-bit settings quantization.
- Goal Achieved: The 1.1B model proved capable of running efficiently with a memory footprint under 1GB, thanks to int4 quantization. It also exceeded the reasoning performance of other open-source models of similar size benchmarks.
- Automotive Improvement: This confirms that Zonal Controllers are hardware-feasible. An NXP S32G or TI TDA4 usually provides access to 4-8GB of LPDDR4 RAM. A 1GB footprint ensures sufficient space for the OS, network buffers, and real-time control stacks, demonstrating that the Resource Contention challenge can be addressed through appropriate modeling compression.

## 5. Conclusion

The integration of Zonal Architectures, Small Language Models (SLMs), and Neuro-Symbolic AI marks a significant breakthrough in automotive engineering, bringing the industry nearer to the goal of fully software-defined, autonomous vehicles. This research demonstrates that the physical and logical limitations of the automotive edge such as latency, power, and safety can be managed by transitioning from centralized, cloud-based LLMs to decentralized, edge-native models SLMs.

The analysis shows that SLMs such as Phi-3 and TinyLlama, when optimized with quantization and knowledge distillation, fit within the computational limits of modern Zonal Controllers like NXP S32G. This supports a distributed intelligence setup where data is processed semantically at the source, reducing bandwidth overload on the vehicle's backbone. However, the inherently stochastic behavior of these models presents a significant safety concern. The case studies of RoboGuard and SYNAPSE strongly indicate that Neuro-Symbolic AI particularly the use of symbolic logic guardrails is essential to enable safe generative AI vehicles.

Future research should focus on standardizing these symbolic interaction layers within the AUTOSAR framework. Creating a standardized "Safe AI Interface" that specifies how probabilistic models interact with deterministic safety monitors will be crucial for gaining regulatory approval. As 2026 nears, OEMs that can implement this
[12]

"Neuro-Symbolic Zonal Architecture" striking a balance between AI's creative capabilities and the strict safety requirements will have a competitive edge automobile.

## References

[1] Texas Instruments, "How a Zone Architecture Paves the Way to a Fully Software-Defined Vehicle," TI Tech. White Paper, Oct. 2024. [Online]. Available: https://www.ti.com/lit/spry345

[2] NXP Semiconductors, "Zonal Architecture & Software Defined Vehicles," NXP Training Pres., 2024. [Online]. Available:(https://www.nxp.com/docs/en/training-presentation/TP-TD24-EUF-AUT-T4753.pdf)

[3] Ravichandran, Z., Robey, A., Kumar, V., Pappas, G. J., & Hassani, H. (2025). Safety Guardrails for LLM-Enabled Robots. *ArXiv*. https://arxiv.org/abs/2503.07885

[4] Z. Lu et al., "Small Language Models: Survey, Benchmark, and Case Studies," ACL Anthology, 2025. [Online]. Available: https://aclanthology.org/2025.acl-long.718.pdf

[5] Castagnone, A., & Nitti, G. (Jan 2026). A Neuro-Symbolic Framework for Ensuring Deterministic Reliability in AI-Assisted Structural Engineering: The SYNAPSE Architecture. Buildings, 16(3), 534. https://doi.org/10.3390/buildings16030534

[6] Abdin, M., Aneja, J., Awadalla, H., Awadallah, A., Awan, A. A., Bach, N., Bahree, A., Bakhtiari, A., Bao, J., Behl, H., Benhaim, A., Bilenko, M., Bjorck, J., Bubeck, S., Cai, M., Cai, Q., Chaudhary, V., Chen, D., Chen, D., . . . Zhou, X. (2024). Phi-3 Technical Report: A Highly Capable Language Model Locally on Your Phone. ArXiv. https://arxiv.org/abs/2404.14219

[7] Zhang, P., Zeng, G., Wang, T., & Lu, W. (2024). TinyLlama: An Open-Source Small Language Model. ArXiv. https://arxiv.org/abs/2401.02385

[8] Texas Instruments, "TSN in Automotive Zone Architectures," TI White Paper, Oct. 2025. [Online]. Available: https://www.ti.com/lit/wp/spry352/spry352.pdf

[9] Z. Ravichandran et al., "RoboGuard: Safety Guardrails for LLM-Enabled Robots," arXiv preprint arXiv:2503.07885, 2025. [Online]. Available: https://arxiv.org/pdf/2503.07885v1.pdf

[10] Dong, Y., Mu, R., Zhang, Y., Sun, S., Zhang, T., Wu, C., Jin, G., Qi, Y., Hu, J., Meng, J., Bensalem, S., & Huang, X. (2025). Safeguarding large language models: A survey. Artificial Intelligence Review, 58(12), 382. https://doi.org/10.1007/s10462-025-11389-2

[11] Microsoft, "AI-Powered Success: Japan Airlines and Phi-4," Microsoft Cloud Blog, July 2025. [Online]. Available: https://www.microsoft.com/en-us/microsoft-cloud/blog/2025/07/24/ai-powered-success-with-1000-stories-of-customer-transformation-and-innovation/.