# Real-Time AI Inference at the Edge for Self-Driving Cars

Self-driving cars represent one of the most challenging applications of edge computing and AI. Modern autonomous vehicles must process over 1 GB of sensor data per second from high-resolution cameras, LiDAR, and radar units to make split-second decisions that ensure passenger safety.

This presentation explores how real-time AI inference at the edge enables autonomous vehicles to function safely and efficiently, examining the computational challenges, hardware innovations, and optimization techniques that make this possible.

By: **Murali Krishna Reddy Mandalapu**

# The Data Challenge of Autonomous Vehicles

**1** Massive Data Volume

Autonomous vehicles generate between 1.4 TB to 19 TB of raw data per hour from multiple high-resolution cameras (30-60 FPS at 1920×1080), LiDAR (100,000 to 4.5 million points per frame), and radar systems operating at 24-77 GHz.

**2** Strict Latency Requirements

The complete perception-decision-action pipeline must execute within 100-300 milliseconds for collision avoidance at highway speeds. Each 10 ms of processing delay translates to approximately 0.3-0.5 meters of extra stopping distance.

**3** Environmental Variability

Computational load can vary by up to 480% between minimal-complexity scenarios (open highways) and maximum-complexity environments (dense urban intersections), requiring adaptive computing architectures.

# Evolution of Edge Computing Hardware

**First Generation: Repurposed Consumer GPUs**

**1**

Early autonomous prototypes used adapted consumer GPUs delivering 8-12 TOPS with significant limitations: high power consumption (250-300W), thermal challenges requiring liquid cooling, and data transfer bottlenecks consuming up to 67% of processing time.

**2**

**Second Generation: Automotive-Grade Accelerators**

Purpose-built automotive processors improved energy efficiency (2-4 TOPS per watt), reduced memory traffic by 40-60% through pruning and compression, and supported reduced-precision computing for 3-4× throughput improvement.

**Current Generation: Heterogeneous Computing**

**3**

Today's platforms combine specialized processors optimized for specific workloads, achieving 10-50× efficiency improvements. They implement dynamic voltage and frequency scaling, reducing power by 30-45% during less intensive scenarios.
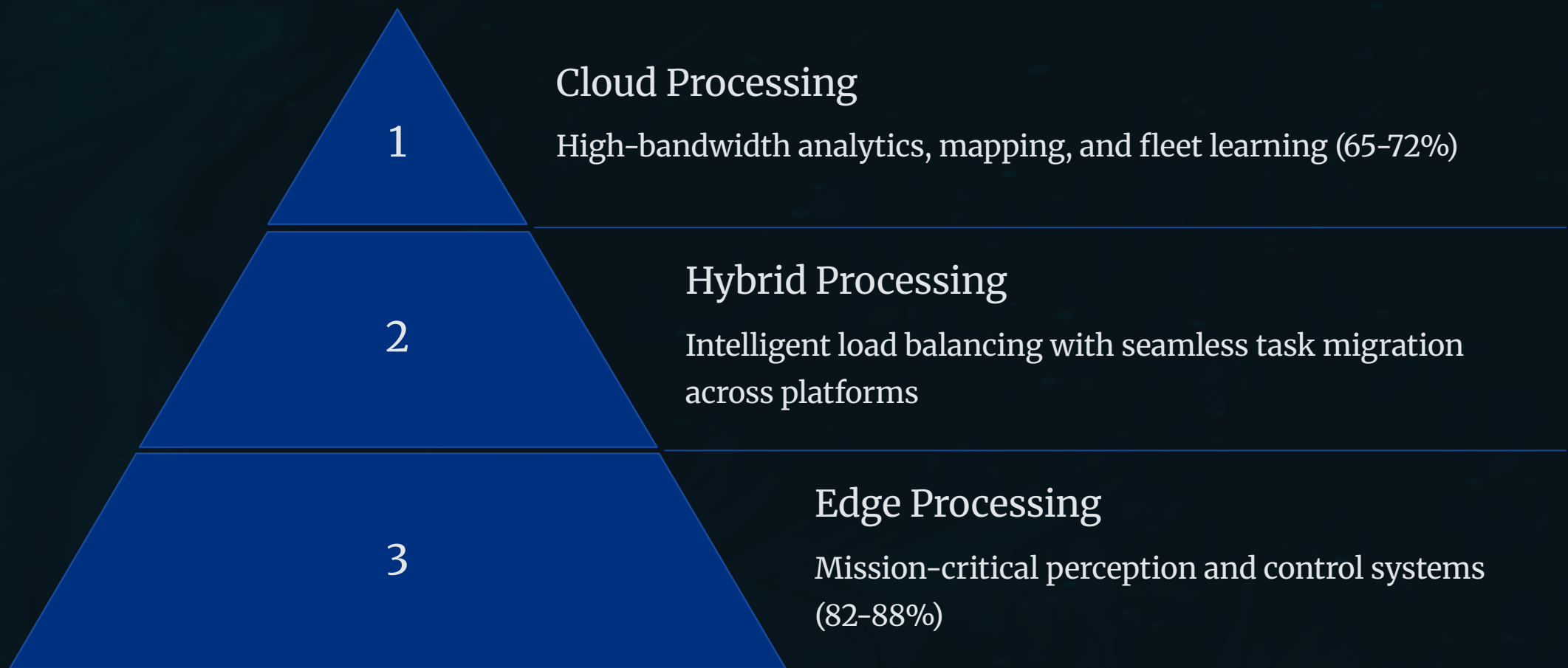
# Edge vs. Cloud Computing Tradeoff

## Edge Computing Advantages

- Near-instantaneous processing (5-15ms latency vs. 50-500ms for cloud)

- Maintains functionality during connectivity interruptions (23-38% of routes)

- Reduced security vulnerabilities (5-7 attack vectors vs. 14-18 for cloud)

- Better privacy protection by keeping sensitive data within vehicle boundaries

## Cloud Computing Advantages

- 2-3 orders of magnitude greater computational throughput

- Enables more sophisticated algorithms with higher accuracy

- Externalizes power consumption (800-1500W for edge processing)

- Ideal for non-safety-critical tasks like HD map generation

# Hybrid Approach: The Current Consensus

**Cloud Processing**

1

High-bandwidth analytics, mapping, and fleet learning (65-72%)

**Hybrid Processing**

2

Intelligent load balancing with seamless task migration across platforms

**Edge Processing**

3

Mission-critical perception and control systems (82-88%)

The industry has converged on a hybrid architecture that strategically allocates computational workloads between onboard systems and cloud infrastructure. This sophisticated approach delivers 99.98% functional reliability while reducing vehicle computational requirements by 30-45% compared to pure edge solutions.

By leveraging both paradigms' inherent strengths—the immediacy of edge processing with the scalability of cloud computing—this hybrid framework creates an optimal balance of safety, efficiency, and capabilities while effectively neutralizing the limitations of either approach used in isolation.

# Object Detection and Classification at the Edge

## Performance Metrics

State-of-the-art models achieve 82-87% mean Average Precision (mAP) while operating within strict latency constraints of 30-50ms per frame on automotive-grade processors.

## Multi-Stage Processing

Initial region proposal networks operate at 20-30Hz, followed by more intensive classification networks that process only identified regions of interest, reducing computational requirements by 65-75%.

## Optimization Techniques

Quantization to 8-bit integer precision reduces memory requirements by 73-76% and inference time by 2.5-3.4× with accuracy degradation of only 1.2-1.8% compared to full-precision models.

# Lane Detection and Trajectory Planning

### 1  Lane Detection

Consumes 8–12% of perception budget, processing camera feeds at 1280×720 to 1920×1080 resolution to extract lane markings with ±5–8cm accuracy at distances up to 80 meters.
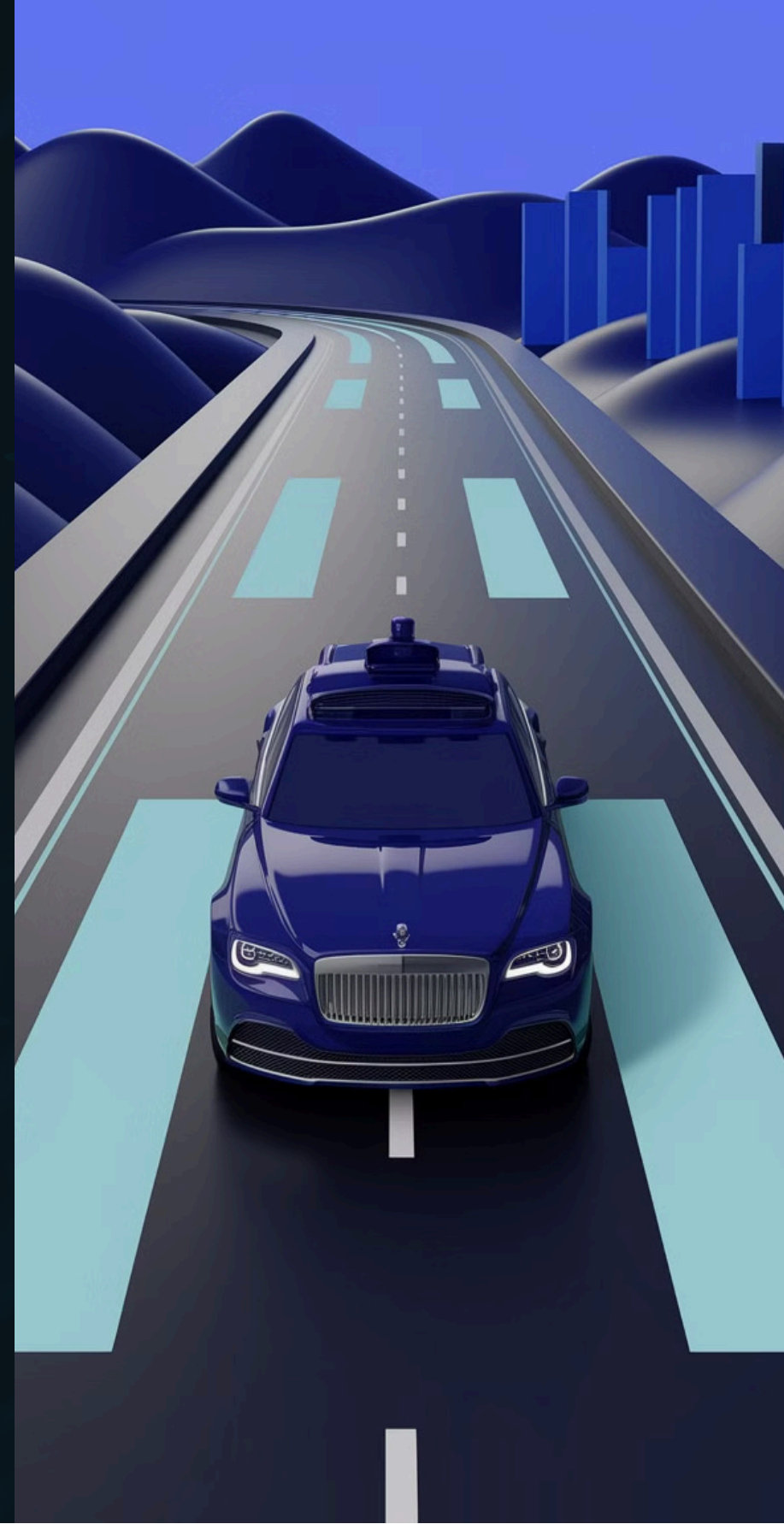
### 2  Multi–Modal Fusion

Combines RGB camera data with LiDAR reflectivity and radar returns to maintain accuracy during adverse weather when visual data quality degrades by 45–60%.

### 3  Trajectory Planning

Evaluates 1,500–3,000 candidate trajectories every 100ms (15–25% of computational budget), optimizing for safety, comfort, and efficiency within strict time constraints.

# Model Optimization Techniques

## Quantization

Converts high-precision 32-bit floating-point to efficient 8-bit integer representation, delivering 4× weight compression and 2× activation compression with minimal accuracy loss (0.5-1.2%). Critical for edge deployment, 8-bit operations consume 9× less energy than floating-point calculations.
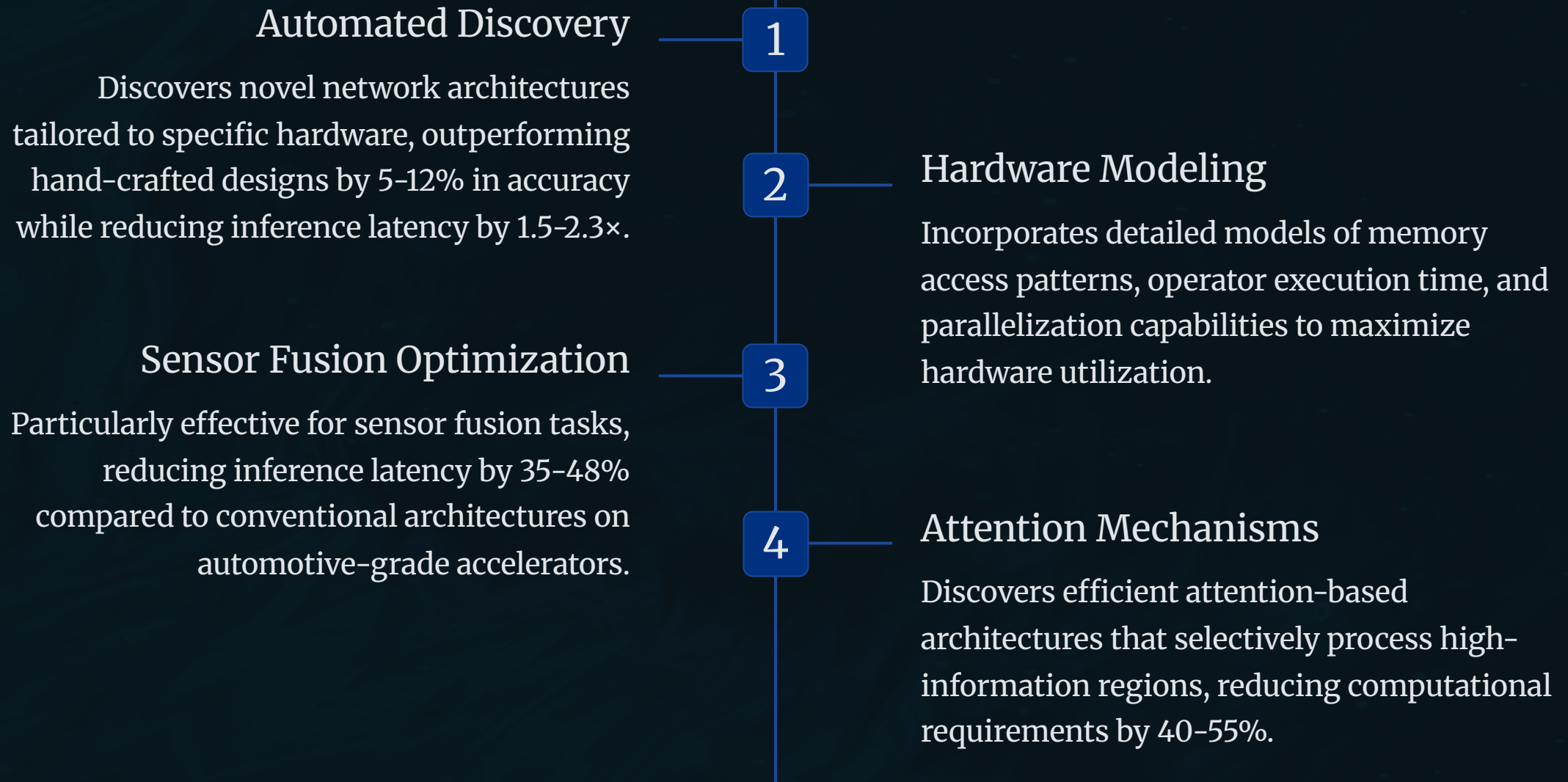
## Pruning

Systematically eliminates redundant neural network parameters (30-70% typically contribute negligibly to performance). Structured pruning techniques yield 2.7-3.8× computational efficiency gains while maintaining model integrity, with accuracy drops under 2% after fine-tuning.

## Knowledge Distillation

Leverages larger "teacher" models to guide the training of compact "student" networks. This technique enables dramatically smaller models (5-8× fewer parameters) to capture the essential capabilities of larger architectures while sacrificing only 2-3% accuracy, ideal for resource-constrained edge devices.

# Hardware-Aware Neural Architecture Search

**Automated Discovery** — **1**

Discovers novel network architectures tailored to specific hardware, outperforming hand-crafted designs by 5-12% in accuracy while reducing inference latency by 1.5-2.3×.

**2** — **Hardware Modeling**

Incorporates detailed models of memory access patterns, operator execution time, and parallelization capabilities to maximize hardware utilization.

**Sensor Fusion Optimization** — **3**

Particularly effective for sensor fusion tasks, reducing inference latency by 35-48% compared to conventional architectures on automotive-grade accelerators.

**4** — **Attention Mechanisms**

Discovers efficient attention-based architectures that selectively process high-information regions, reducing computational requirements by 40-55%.

# Emerging Trends: Neuromorphic Computing

| 1 | ### Event-Driven Processing |
|---|---|
| | Reduces power consumption by 90–95% compared to frame-based approaches by allocating resources only when significant changes occur. |

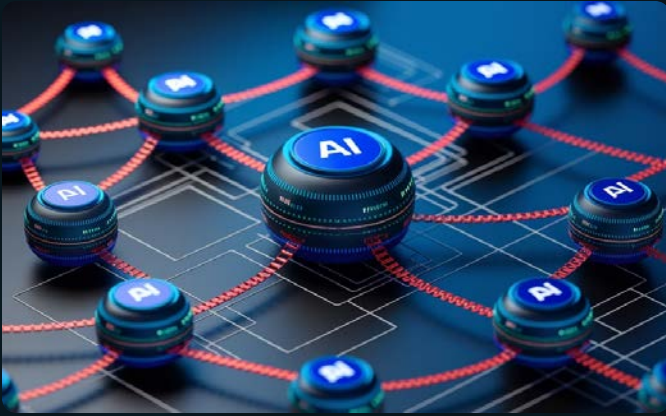| 2 | ### Temporal Advantages |
|---|---|
| | Microsecond-scale temporal resolution (1–10µs) improves detection latency by 20–45ms for high-speed objects compared to conventional vision systems. |

| 3 | ### Lighting Adaptability |
|---|---|
| | Maintains consistent detection across illumination ranges from 0.1 lux to 100,000 lux, addressing limitations with high dynamic range environments. |

Neuromorphic computing represents a fundamental departure from conventional architectures, drawing inspiration from biological neural systems. These systems integrate processing and memory in artificial neuron and synapse structures that mimic their biological counterparts.

Prototype neuromorphic processing units can process 50–100 million events per second while consuming only 100–300mW of power, representing a two orders of magnitude improvement in efficiency compared to GPU-based solutions.

# The Future: Distributed AI and Continuous Learning



### Distributed AI Architectures

Maintain 85-92% of critical functionality even with failures in 30% of processing nodes. Reduce latency by 35-47% for complex perception tasks through better parallelization and reduced data movement.



### Continuous Learning Systems

Improve detection accuracy by 15-25% in novel environments not represented in initial training data. Employ safety-aware incremental updates that limit parameter changes to 0.5-2% per cycle.



### Federated Learning

Vehicles contribute to collective intelligence while transmitting only 0.1-0.5% of data required for centralized approaches. Synchronization every 100-500 kilometers provides 80-85% of continuous connectivity benefits.

Thank you