

Agentic AI Scene Parsing Master Class Q&A with Will Gray-Roncal

1) Question: Can generative AI models (for example, diffusion models and Vision Transformers) do object detection better, or are traditional computer-vision models like convolutional neural networks (CNNs) and You-Only-Look-Once (YOLO) still preferred for these jobs?

Answer: For real-time object detection with a known label set, dedicated detectors (for example, YOLO or DEtection TRansformers, DETR) are usually faster and more accurate at the same cost. Foundation models (for example, Vision Transformers or open-vocabulary detectors) shine when the label space is open-ended or changes over time. Practical rule: if classes are stable and you have labeled data, ship a dedicated detector; if classes drift or are long-tail, consider an open-vocabulary detector and optionally distill it into a smaller runtime model.

2) Question: Did image recognition become easier to handle after the transformer architecture became mainstream? If so, how?

Answer: The largest gains came from large-scale pretraining and self-supervision (for example, masked-image pretraining and contrastive language–image pretraining), not only the Vision Transformer block itself. These approaches produce features that transfer well across tasks. If compute is tight, strong convolutional networks remain competitive; if you need cross-modal reasoning or long-range context, transformer backbones simplify the overall system.

3) Question: How do these empirical methods align with the principles of Responsible and Ethical AI?

Answer: Treat ethics as engineering. Encode intended and disallowed uses in tests; record data sources and consent; evaluate performance on demographic and scenario slices; red-team before deployment; monitor for drift; and keep a rollback plan. Automate these practices and make them visible on dashboards so they influence daily decisions.

4) Question: Is this how the brain works—do artificial networks understand and analyze data like the brain and differentiate between colors?

Answer: No. These systems are inspired by biology but are not faithful replicas. Convolutions echo local receptive fields and attention resembles dynamic routing, but artificial networks operate on numbers, not spikes and neurotransmitters. They are useful abstractions, not working brain models.

5) Question: In the future, when we have robots or digital agents without human-like faces, how will “facial recognition” evolve to identify or authenticate them?

Answer: Shift identity from faces to cryptographic attestation, device certificates, and behavioral fingerprints (for example, motion patterns or radio-frequency signatures). Use visual perception only as a secondary check; make strong identity a protocol property, not a camera trick.

6) Question: What does “transformer” mean here? Is it simply transferring data, or something else?

Answer: “Transformer” refers to a neural-network architecture that replaces recurrence with self-attention, allowing each token or image patch to weigh information from all others in one step. The name reflects how it transforms sequences via attention-based mixing.

7) Question: For the mechanics of image recognition and computer vision, how many kinds of image analysis are common, and how do we combine them into a reliable result?

Answer: Common tasks include detection, segmentation (semantic, instance, and panoptic), keypoints and pose, depth or surface normals, optical character recognition, retrieval, captioning, and visual question answering. Combine them with a shared encoder and task-specific heads or via late fusion with validators. Reliability improves when you surface disagreements and route edge cases to a human reviewer.

8) Question: How do you envision agentic AI shaping the future of Mixed Reality (MR) and artificial-intelligence integration, especially in computer vision?

Answer: Mixed Reality needs grounded perception, spatial memory, and tool use. An agent can maintain a scene graph, reason about tasks, and call planners or checklists. Start with assistive workflows (for example, inspection or assembly) and measure task-completion time and human handoffs.

9) Question: Techniques like CAPTCHA have been used to aid image recognition. Can you explain that?

Answer: Early systems like reCAPTCHA generated large-scale human labels that helped train recognition models. Today, prefer purpose-built datasets or self-supervision with explicit consent; avoid harvesting labels without clear permission.

10) Question: Which is more energy-efficient—biological neurons or artificial neurons? Current power usage suggests artificial neurons are less efficient.

Answer: Biology is far more energy-efficient—often orders of magnitude better per operation. We narrow the gap with sparsity, low-precision arithmetic, specialized accelerators, and on-device inference to reduce data movement.

11) Question: Many information-technology service providers claim multi-agent solutions (for example, support agents). Are these suites built on top of existing, market-available AI technologies?

Answer: Usually yes. Most offerings orchestrate large language models, retrieval systems, and existing tools behind a workflow layer. Ask vendors about their evaluation harness, rollback plan, cost controls, data isolation, and measured uplift on your tasks.

12) Question: Where do you invest the most when building agentic AI? How do you decide between agentic AI, traditional machine learning, or conventional automation?

Answer: Invest in task decomposition, clean tool interfaces, evaluation, observability, and feedback loops. Decision guide: deterministic workflow → robotic process automation; stable input-output mapping with data → traditional machine learning; open-ended goals with tool use → an agentic system. Define service-level objectives first in all cases.

13) Question: Is debugging harder in agentic AI than in traditional AI?

Answer: Yes. Non-deterministic planning and tool chains create more failure points. Mitigate with deterministic fixtures, frozen random seeds, rich trace logs of plans and tool calls, and counterfactual replays so you can reproduce and fix issues.

14) Question: For those with a microservices background, how should we think about multi-agent systems? Are agents similar to microservices, or fundamentally different?

Answer: Keep tools as microservices with clear contracts. The agent is a planner that decides which service to call and in what order to reach a goal. That adds flexibility but requires strong observability at both layers.

15) Question: How do today's topics connect—from image parsing to neural networks to effective prompts? Is there a theme?

Answer: The through-line is grounding → reasoning → acting. Vision grounds the world, language models reason, and agents act via tools. Prompts, memory, and evaluation tie these into one system.

16) Question: Should I use Pinecone as the default vector database, or is Chroma with the OpenAI API acceptable?

Answer: Choose based on operations. Managed vector stores such as Pinecone or Weaviate help at scale and enable hybrid search. Chroma is great for local prototypes. PostgreSQL with the pgvector extension is a strong default until latency or scale demands a move.

17) Question: Are there tools to measure an agent's confidence level, or should agents self-assess?

Answer: Use task metrics (success rate and time to resolution), model-uncertainty proxies (token probabilities or entropy), agreement across multiple runs, and external validators that check units, schemas, or constraints. Do not rely on a single “confidence score.”

18) Question: How do we evaluate an agent's confidence and correctness?

Answer: Require evidence. For retrieval-augmented systems, link answers to sources and check faithfulness. For tool use, keep contract tests. For math or code, run executable tests. A practical pattern is proposer → critic → verifier, with logs explaining each decision.

19) Question: Are these agent characteristics applicable outside of image recognition?

Answer: Yes. Planning, memory, and tool use generalize to data engineering, security triage, policy analysis, education technology, and more. Vision is simply a concrete teaching anchor.

20) Question: In what ways can interactive agents—such as social robots and AI-driven conversational systems—help children with autism spectrum disorder build communication and social skills? Please share a few examples.

Answer: Useful functions include turn-taking games, emotion coaching, and visual schedules with gentle prompts. Keep clinicians or parents in the loop, protect privacy, and

avoid autonomous diagnosis. Examples include social-story companions and guided-practice apps that track progress.

21) Question: How will agentic AI shape the future of Mixed Reality and artificial-intelligence integration, given our focus on computer vision?

Answer: Strong fit: the agent maintains spatial context, converses with the user, and calls planners or measurement tools. Main risks are latency, tracking drift, and cognitive overload—solve with lightweight on-device models and graceful cloud fallback.

22) Question: Are there design patterns for multi-agent planning?

Answer: Yes. Common patterns include single orchestrator (a router-planner), proposer-critic (writer-editor), supervisor with specialists (hub-and-spoke), marketplace or blackboard (publish-subscribe with bidding), and self-play for exploration. Choose based on decomposition clarity, redundancy needs, and cost.

23) Question: We are moving into using large language models wrapped with agent modes (prompting, retrieval-augmented generation, and application-programming interface data sourcing) and need near real-time evaluation. Which evaluation tools work best for monitoring and benchmarking?

Answer: Start with an in-house harness that replays logged scenarios, perturbs inputs, and tracks latency, tool-error rate, fallback rate, and success. Add switchback A/B testing when you have traffic. Simple, visible dashboards beat heavy frameworks at the beginning.

24) Question: How should I choose large language models when the task involves image recognition?

Answer: Language models alone are poor at pixel-level perception. Use detectors or segmenters for images; use a vision-language model to describe or ground findings; then apply a large language model for reasoning. Keep measurable checks at each step. For some scene-parsing tasks, modern multimodal models can perform well—evaluate empirically.

25) Question: How is agentic AI different from robotic process automation? I've heard that agents handle non-deterministic problem solving, unlike robotic process automation or microservices.

Answer: That intuition is right. Robotic process automation is scripted determinism. Agents pursue goals under uncertainty by planning and choosing tools. Still embed deterministic subroutines and guardrails inside the agent pipeline.

26) Question: What platform or tools were used to build this demo?

Answer: A typical stack is a Python or Node.js orchestrator, lightweight tool adapters over HTTP, and a React web interface with a trace panel so we can replay plans and tool calls.

27) Question: How did you build this user interface? Was it programmed in Python?

Answer: The user interface was built with a web framework (for example, React) and component libraries. Python handled back-end orchestration and evaluation; the front end communicates with it over web endpoints.

28) Question: What evaluation tools are used in this demo to assess confidence and accuracy?

Answer: Offline we use scenario sets with known answers, contract tests, hallucination checks, and cost and latency profiling. Online we track success rate, first-pass yield, fallback percentage, human-handoff rate, and alerts for tool failures or looping. We also monitor for bias and other responsible-AI concerns.

Thanks for the great questions!