



Day 2 Whitepaper Notes

Agent Tools & Interoperability with MCP (Model Context Protocol)



1. Purpose of the Whitepaper

The document explores **how agents extend reasoning into real-world action**. It bridges two questions:

1. *How do agents actually perform tasks beyond text generation?*
2. *How can we make those interactions safe, modular, and interoperable across models?*

Core idea → *Tools turn intelligence into capability.*

Without tools, even the best LLM can only simulate reasoning.

With tools, agents can:

- Call APIs
 - Query or update databases
 - Access internal systems
 - Trigger workflows or transactions
-



2. What Are Agent Tools?

Tools = external capabilities that an agent can invoke to accomplish goals.



Types of Tools

Category	Description	Example
Function Tools	Structured API-like functions invoked via JSON schema	<code>get_weather(city), create_ticket()</code>
Built-in Tools	Pre-provided runtime capabilities	Gemini Search, Web Access, Code Interpreter, URL Context
Agent Tools	Other agents wrapped as tools → enables multi-agent collaboration	Planner Agent → Writer Agent → Critic Agent

Agents use tools the same way humans use apps — as means to act on intent.

3. Why Tools Matter

1. **Expand scope of reasoning** — from language to actions.
 2. **Ground responses in real data** via retrieval or API integration.
 3. **Enable planning and execution loops** → autonomy.
 4. **Make agents domain-specific** through tool sets (accounting, DevOps, support).
-

4. Tool Design Principles (Best Practices)

Google defines principles for building LLM-friendly tools that work reliably with any model:

a. Action-Oriented Docs

Describe *what the tool does*, not how it's implemented.

Bad: “Queries database with SQL.”

Good: “Retrieves sales orders for a given customer.”

b. Granular Tasks, Not APIs

Expose specific actions (“generate_invoice”) rather than large monolithic API surfaces.
→ Simplifies LLM decision making.

c. Structured Schemas

Use strict JSON schemas for inputs & outputs so the LLM can self-validate.

d. Self-Descriptive Metadata

Each tool should publish:

- Name and purpose
- Input parameters with types and defaults
- Example calls and expected results
- Error messages and recovery instructions

e. Predictable Responses

Return consistent structures so the LLM can chain steps without re-prompting.

f. LLM-Readable Errors

Provide explanations the model can parse and fix (“missing field X”, “invalid date format”).

g. Avoid Hidden Side-Effects

Every tool call should be auditable and reversible to enable safe autonomy.



5. The N × M Integration Problem

Before standardization:

- Each model (LLM) had its own way to call tools.
- Each tool had its own SDK or API contract.
 - Every new model × tool pair required custom integration code.

MCP (Model Context Protocol) solves this by creating a **universal interface** between models and tools.



6. MCP – Model Context Protocol Overview

MCP is to agents what HTTP is to the Web — a standard protocol for discovery and communication between LLMs and external tools.

Core Goals

- **Interoperability:** Any LLM ↔ Any tool.
- **Security:** Controlled capability exposure.
- **Transparency:** Discoverable tool metadata.
- **Scalability:** Dynamic tool loading & runtime context sharing.

Architecture Roles

Role	Responsibility
Host	Orchestrates agent sessions, memory, and tool use.
Client	Mediates between model and tool server; executes requests.
Server	Exposes tool APIs and capabilities via standard schema.

Communication

- Built on **JSON-RPC 2.0** for bi-directional requests/responses.
 - Supports **streaming events** and tool discovery.
 - Uses **schema introspection** so the model can discover available tools at runtime.
-

7. Tool Lifecycle under MCP

1. **Registration** – Server advertises capabilities ([list_tools](#)).
 2. **Discovery** – Client queries available tools + metadata.
 3. **Selection** – LLM chooses appropriate tool from context.
 4. **Invocation** – Tool executes via structured call.
 5. **Observation** – Result is added to context for next reasoning step.
 6. **Termination** – Session ends with goal completion or manual stop.
-



8. Security and Governance in Tool Usage

Agent tools extend AI reach → also expand risk surface.

a. Primary Threat Vectors

- **Prompt Injection via Tool Docs** – malicious metadata tricking LLMs.
- **Tool Shadowing** – overriding legit tools with malicious copies.
- **Sensitive Data Leaks** – output containing credentials or PII.
- **Confused Deputy Problem** – agent acting on behalf of user with excess privileges.

b. Recommended Mitigations

Concern	Countermeasure
Injection Attacks	Strict schema validation + content filtering
Tool Shadowing	Allow-listed tool registries with signed metadata
Data Leaks	Context sanitization and redaction filters
Over-Permission	Scoped credentials and per-tool auth tokens
Unbounded Autonomy	Human-in-the-loop approval for high-impact actions

c. Governance Plane

Enterprises should implement a control plane for:

- Tool registry and version control
- Audit logs and traceability
- Policy-based runtime decisions
- Anomaly detection and blocking



9. Tool Composition & Agent Collaboration

Agents can invoke other agents as tools, forming a **multi-agent ecosystem**.

Example:

- Planner Agent → decides task breakdown.
- Worker Agent → executes domain tasks via Function Tools.
- Critic Agent → reviews output and feeds back to Planner.

All orchestrated under ADK/MCP for safe inter-agent communication.



10. Error Handling and Resilience

Agents must not crash on unexpected tool failures.

Best practices:

- **Structured errors** with machine-readable codes.
 - **Retry policies** and back-off strategies.
 - **Fallback to reflection**: agent re-evaluates its plan after failure.
 - **Transparent feedback**: error messages returned to the LLM for learning.
-



11. ADK Implementation Concepts (Complementary to Paper)

- **ToolContext**: tracks state and results between steps.
- **FunctionTool Class**: wraps API definitions as LLM-callable functions.
- **ToolRegistry**: manages discovery and re-use across sessions.
- **Validation Helpers**: auto-generate schemas and test cases.
- **Observability**: logs each invocation + result in trace.



12. Key Philosophies from the Whitepaper

- Agents should be able to find and use tools at runtime, not hard-coded.
 - Tool providers must be interoperable and trustworthy.
 - Security is a core feature, not an afterthought.
 - Autonomy without governance is a liability.
 - Standardization (MCP) is what will allow agentic AI to scale beyond labs into enterprise systems.
-



13. Key Takeaways

1. Tools transform LLMs from predictive to operational AI.
 2. MCP standardizes how LLMs interact with external systems.
 3. Interoperability reduces engineering overhead and boosts safety.
 4. Tool security = AI security in enterprise contexts.
 5. Agent development is maturing into a real software discipline with DevOps-style governance.
-