

Model Context Protocol (MCP)–Powered Distributed Web Crawling System: A Framework for Guided, Agentic Crawling at Scale

Executive Summary

The emergence of AI agents as primary computing paradigms has fundamentally altered how organizations access and process web-scale information. Traditional web crawlers—built on stateless, link-following heuristics and centralized schedulers—are increasingly misaligned with modern agentic architectures. This document presents **MCP-Guided Web Crawling (MGWC)**: a next-generation framework for distributed, autonomous web exploration where **Guide Bot agents** (powered by frontier LLMs) make intelligent, context-aware decisions about what to crawl, how to crawl it, and how to collaborate with peer agents.

The framework is built on the **Model Context Protocol (MCP)**, an open standard enabling AI agents to securely access external tools, data, and policies. By adopting MCP as the lingua franca for web crawling infrastructure, organizations including OpenAI, Perplexity, Google, Microsoft, and others can build **interoperable, scalable crawling systems** where:

- **Guide Bot agents** reason about link importance, content relevance, and safety using real-time context.
- **Shared MCP servers** expose crawling primitives (frontier management, URL deduplication, reputation scoring, robots.txt caching) as standardized tools.
- **Distributed agent collectives** coordinate autonomously without centralized bottlenecks, enabling efficient coverage of the modern web.
- **Plug-and-play integration** with existing enterprise systems (Search Index, Safety APIs, Knowledge Graphs, LLM training pipelines).

This document is intended as both a **conceptual reference** and a **practical implementation guide** for engineering teams building production crawling infrastructure. It details:

1. **Architecture & Design Principles** – Why MCP is the right abstraction for agentic crawling.
2. **Real-World Use Cases** – How OpenAI, Google, Microsoft, Perplexity, and others can apply MGWC today.
3. **MCP Server Specification** – Detailed schemas and tools for Guide Bot-compatible crawling services.
4. **Agent Coordination Patterns** – Multi-agent collaboration, conflict resolution, and resilience.
5. **Implementation Roadmap** – Phased deployment strategies for existing crawlers.
6. **Production Considerations** – Safety, compliance, observability, and performance at scale.

1. Introduction: The Case for MCP-Guided Web Crawling

1.1 The Problem with Traditional Crawling in the AI Era

For decades, web crawlers have operated on a simple principle: fetch every reachable URL, parse it, extract links, and repeat. This model worked well when:

- Search engines needed exhaustive coverage.
- Crawling was I/O-bound and could ignore content semantics.
- A single centralized scheduler could manage global state.
- Web scale was measured in billions of pages, not trillions of entities.

Today's landscape is radically different:

1. Agentic Workloads Demand Semantic Intelligence

Modern AI agents (ChatGPT, Claude, Gemini, Copilot) don't just index—they reason, understand, and make autonomous decisions. When an agent is tasked with "find the latest research on quantum error correction" or "compare pricing for enterprise SaaS across 5 vendors," the crawl must be:

- **Semantically targeted:** Skip spam, marketing fluff, and outdated content. Prioritize authoritative sources.
- **Real-time adaptive:** If a new source emerges mid-crawl, the agent should know about it within minutes.
- **Trust-aware:** Understand which sources are authoritative, which are promotional, which are compromised.

Traditional crawlers optimize for *coverage*. Agentic crawlers must optimize for *relevance* and *reliability*.

2. Distributed Agentic Systems Need Coordination

A single LLM-powered chatbot is inefficient. Future systems will be **agent collectives**: hundreds of specialized agents (one for financial news, one for academic research, one for product catalogs, one for infrastructure monitoring) crawling concurrently.

Without coordination:

- 50 agents crawl the same Tesla earnings report independently, wasting bandwidth.
- One agent discovers Tesla's new supplier network; others don't learn it.
- Agent A crawls tesla.com at 100 req/s while Agent B crawls at 50 req/s; the host sees 150 req/s (polite crawling violated).

3. Modern Web Complexity Requires Specialized Tools

Today's web is not simple HTML:

- JavaScript-rendered content (React, Vue, dynamic SPAs).
- Bot detection and anti-scraping (Cloudflare, DataDome, Imperva).
- Federated authentication (OAuth, SAML, SSO).
- Multilingual, multimedia content (video, audio, PDFs, embedded APIs).
- Constantly shifting rate limits and access policies.

Monolithic crawlers cannot adapt. Specialized MCP-based tools (Puppeteer for JS rendering, credential managers for auth, ML-based bot detection evasion, multilingual parsers) can be composed on demand.

4. Fragmentation Across Organizations

OpenAI, Google, Microsoft, Perplexity, and others each build custom crawling stacks. Each reinvents:

- URL deduplication.
- Reputation scoring.
- Safety checks.
- Robots.txt parsing.
- Crawl scheduling.

This is **massively wasteful**. A shared, open standard (MCP) could allow these organizations to:

- Contribute shared tools (e.g., "Google's safe browsing MCP server").
- Interoperate on crawling campaigns.
- Pool threat intelligence.

1.2 Why MCP is the Answer

The **Model Context Protocol** was designed to solve precisely this problem: enabling AI agents to securely access external systems without requiring custom integrations for every new data source or capability.

MCP provides:

Standardization: A universal schema (JSON-RPC 2.0, tools, resources, prompts) means any MCP-aware system can discover and use any MCP server.

Composability: Rather than monolithic crawlers, crawling becomes a composition of MCP servers:

- One server handles URL frontier management.
- Another handles reputation scoring.
- Another caches robots.txt.
- Another integrates with SafeBrowsing APIs.
- A Guide Bot agent (the client) orchestrates all of them.

Scalability: MCP servers can be run locally (fast, low latency) or remotely (shared infrastructure). They can be replicated, sharded, or cached as needed. The protocol itself is stateless and lightweight.

Interoperability: If Anthropic, OpenAI, and Google all implement MCP-based crawling tools, their agents can cross-call each other's services. This enables **federated crawling networks** where organizations share infrastructure without surrendering autonomy.

Security: MCP includes built-in authentication, authorization, and audit trails. An enterprise can safely expose crawling tools to AI agents (even untrusted agents) knowing that access is controlled and logged.

1.3 This Document's Scope

This document outlines a **complete, production-ready architecture** for MCP-Guided Web Crawling (MGWC). It is aimed at:

- **Engineering leads** at OpenAI, Google, Microsoft, Perplexity, and other AI organizations deciding how to architect next-generation crawling systems.
- **MCP server developers** building crawling-specific tools and integrations.
- **Agent builders** designing autonomous systems that crawl intelligently.
- **Security and compliance teams** ensuring crawling respects legal, ethical, and technical boundaries.

The document will cover:

1. **Architecture and design patterns** for agentic crawling.
2. **Real-world use cases** showing how major organizations can benefit.
3. **Detailed MCP server specifications** for crawling-critical services.
4. **Multi-agent coordination protocols** for large-scale deployments.
5. **Implementation roadmaps** for migrating existing crawlers.
6. **Production hardening** (safety, compliance, monitoring).

2. Core Architecture: MCP-Guided Web Crawling (MGWC)

2.1 System Components

An MGWC system consists of four primary layers, all communicating via MCP:

Layer 1: Web Crawler Agents (Clients)

A web crawler is no longer a monolithic binary. Instead, it becomes an **LLM-powered agent** (a "Guide Bot") that reasons about what to crawl next.

Characteristics:

- Small, lightweight process (can run in a container, on a laptop, or in the cloud).
- Consumes HTTP/HTML pages and extracts links (traditional crawling functions).
- For each discovered link, queries MCP servers to decide: "Should I crawl this?"
- Maintains minimal local state (current page, immediate frontier).
- Reports discoveries and completion to MCP servers.

Example pseudocode:

```
class GuideBotCrawler(Agent):
    def __init__(self, mcp_client, topic="financial_news"):
        self.mcp = mcp_client
        self.topic = topic
        self.local_frontier = []

    def crawl_page(self, url):
        html = self.fetch(url) # Traditional HTTP fetch
        links = self.extract_links(html)
```

```

# Query MCP servers to decide which links to follow
for link in links:
    check_resp = self.mcp.call_tool(
        "frontier.score_url",
        {
            "url": link,
            "context": {
                "source_page": url,
                "topic": self.topic,
                "depth": self.depth(url),
            },
        }
    )

    # Only crawl high-scoring links
    if check_resp["score"] > THRESHOLD:
        self.local_frontier.append(link)

    # Report progress to MCP coordination layer
    self.mcp.call_tool(
        "coordination.report_crawl",
        {
            "crawler_id": self.id,
            "url": url,
            "num_links_extracted": len(links),
            "num_links_queued": len(self.local_frontier),
        }
    )
)

```

Layer 2: MCP Servers (Shared Infrastructure)

MCP servers expose crawling capabilities as standardized tools and resources. Multiple organizations can implement these servers.

Core MCP servers in an MGWC system:

1. **Frontier Server** (mcp-frontier-manager)
 - Tools: score_url, get_next_seeds, report_seen, get_domain_policy
 - Resources: domain_metadata, topic_taxonomy, crawl_schedule
 - Manages what's worth crawling and enforces global priorities.

- 2. Reputation Server** (mcp-reputation-db)
 - Tools: get_domain_score, check_url_safety, report_malicious, get_credibility_signals
 - Resources: domain_trust_scores, threat_intelligence_feeds
 - Maintains and shares trust information about domains and URLs.
- 3. Coordination Server** (mcp-crawl-coordinator)
 - Tools: register_crawler, report_crawl_status, request_work_assignment, sync_seen_urls
 - Resources: global_crawl_state, crawler_registry, active_campaigns
 - Orchestrates multi-crawler campaigns, avoids duplication.
- 4. Caching Server** (mcp-crawl-cache)
 - Tools: cache_robots_txt, cache_content_signature, lookup_duplicate, get_etag
 - Resources: robots_cache, content_signatures, url_dedup_registry
 - Reduces redundant fetches and checks.
- 5. Policy Server** (mcp-policy-enforcer)
 - Tools: check_crawl_compliance, get_rate_limits, apply_policy, audit_crawl
 - Resources: robots_txt_rules, crawl_policies, legal_boundaries, compliance_rules
 - Ensures crawling respects robots.txt, legal restrictions, and organizational policies.
- 6. Safety & Detection Server** (mcp-safety-gateway)
 - Tools: scan_for_malware, detect_phishing, check_content_policy, get_safety_score
 - Resources: malware_signatures, phishing_databases, nsfw_classifiers, banned_domains
 - Prevents crawlers from visiting dangerous or prohibited content.

Layer 3: Coordination & State Layer (MCP Clients)

A set of MCP clients orchestrate the crawling campaign at a higher level.

Typical components:

- **Campaign Manager:** An LLM-powered agent that understands crawling goals ("Index all e-commerce sites in Southeast Asia selling EV charging equipment") and delegates to Guide Bot crawlers.
- **Global State Store:** A shared database (Redis, PostgreSQL, or distributed store) that MCP servers read/write. Accessible via MCP for cross-organizational queries.
- **Analytics & Observability:** Tracks crawl metrics (pages indexed, links discovered, errors, safety incidents) and makes this data accessible via MCP resources.

Layer 4: Enterprise Integration

Each organization (OpenAI, Google, Microsoft, Perplexity) integrates MGWC with:

- **Search Index:** New crawled pages → indexed and searchable by agents.
- **LLM Training Pipeline:** High-quality crawled content → fed into training data with proper attribution and licensing compliance.
- **Knowledge Graph:** Discovered entities (people, organizations, products) → added to knowledge graph.
- **Safety & Legal:** Compliance checks, GDPR/CCPA handling, copyright detection.

2.2 Information Flow: An Example Crawl

To concretize, here's how a single Guide Bot crawls:

1. Guide Bot starts, registers with Coordination Server:
→ MCP call: coordination.register_crawler({crawler_id: "bot_001", topic: "academic_research"})
2. Asks for initial seeds:
→ MCP call: frontier.get_next_seeds({topic: "academic_research", num_seeds: 10})
→ Response: [[arxiv.org](#), [scholar.google.com](#), [researchgate.net](#), ...]
3. Fetches first URL, extracts links:
→ HTTP GET [arxiv.org/list/cs.AI/recent](#)
→ Extracts 50 links to papers, authors, institutions
4. For each link, queries frontier and reputation servers:
→ MCP call: frontierscore_url({
url: "[arxiv.org/abs/2401.12345](#)",
context: {source_page: "[arxiv.org/list/cs.AI/recent](#)", topic: "academic_research"}
})
→ Response: {score: 0.92, rationale: "Recent AI paper, high relevance"}
→ MCP call: reputation.get_domain_score({domain: "[arxiv.org](#)"})
→ Response: {trust_score: 0.99, reasons: ["Published by Cornell University", "Long history"]}
5. Only crawls links above score threshold, reports progress:
→ MCP call: coordination.report_crawl_status({
crawler_id: "bot_001",
url_crawled: "[arxiv.org/list/cs.AI/recent](#)",
num_new_urls: 45,
num_skipped: 5,
status: "success"
})
6. Repeat steps 3–5 for next URLs in frontier.
7. During crawl, if another bot reports seeing the same URL:
→ MCP call: coordination.sync_seen_urls({
urls_seen: ["[arxiv.org/abs/2401.12345](#)", ...]
})
→ All bots learn to skip URLs others have seen.

2.3 Design Principles

MGWC is built on several key principles:

1. **Autonomy with Coordination:** Each crawler is autonomous (can make local decisions quickly) but reports to shared MCP servers for global awareness.
2. **Asynchronous & Non-Blocking:** A crawler doesn't wait for MCP responses. It can query one server while fetching from another, using async patterns.
3. **Stateless MCP Servers:** Servers don't track individual crawler state. Crawlers report state; servers maintain global aggregates.
4. **Pluggable Safety & Compliance:** Rather than embedding policy logic in crawlers, policies are MCP tools that can be swapped or extended.
5. **Agent-Native Reasoning:** Guide Bots use LLM reasoning, not hard-coded heuristics. An agent can explain: "I'm skipping this link because the reputation server says the domain is 70% spam, and the topic scorer says it's off-topic for academic research."

-
- 6. Open & Extensible:** New organizations can implement their own MCP servers (e.g., "Microsoft Threat Intelligence Server") and plug into the ecosystem. A Guide Bot automatically discovers and uses new servers.

3. Real-World Use Cases: How Major Organizations Benefit

3.1 OpenAI: Powering ChatGPT's Real-Time Web Access

Current Challenge: ChatGPT's knowledge cutoff and inability to browse the live web limits its utility for time-sensitive queries ("What was the S&P 500 close today?" "What are the latest GPT updates?").

MGWC Solution:

- 1. Deploy financial news crawlers** (Guide Bots specialized for financial domains):
 - MCP frontier server prioritizes [finance.sec.gov](#), [investor.openai.com](#), [bloomberg.com](#), [cnbc.com](#).
 - MCP reputation server identifies which financial news sources are reliable vs. clickbait.
 - Crawlers run every hour, discovering new earnings reports, SEC filings, and analysis.
- 2. Deploy tech news crawlers:**
 - Specialized frontier server for [hackernews.com](#), [techcrunch.com](#), [theverge.com](#), [arxiv.org](#).
 - Real-time prioritization of posts likely to affect AI/LLM development.
- 3. Integration with ChatGPT:**
 - When user asks "What's the latest on OpenAI?", ChatGPT's agent queries the MCP Knowledge Graph server.
 - If content is <1 hour old, served immediately. If >4 hours old, crawlers are triggered to refresh.
 - User gets up-to-date, cited sources.

Benefits:

- Real-time, cited web access without proprietary crawling infrastructure.
- Can outsource crawling to ecosystem partners (others implement MCP frontier/reputation servers).
- Agents reason about freshness vs. compute cost (old data from cache vs. fresh data from new crawl).

3.2 Google: Evolving Search Crawler Infrastructure

Current Challenge: Google's Googlebot crawls billions of pages, but the centralized frontier scheduler is increasingly complex. Incorporating LLM-based ranking and quality signals requires re-architecting the entire system.

MGWC Solution:

- 1. Migrate Googlebot to MCP-Guided architecture:**
 - Replace centralized frontier with an MCP Frontier Server that Googlebot (as a Guide Bot) queries.

- MCP reputation server incorporates Google's PageRank, E-E-A-T signals, and ML-based quality classifiers.
- MCP safety server integrates Google Safe Browsing, malware detection, and spam filtering.

2. Topical crawling specialization:

- Instead of one generic crawler, deploy specialized Guide Bots: news crawler, e-commerce crawler, academic crawler, video index crawler.
- Each subscribes to different MCP frontier servers (news.frontier, ecommerce.frontier, academic.frontier).
- MCP coordination server distributes crawl load across specialties.

3. Multi-organization federation:

- Google publishes its MCP servers as open standards (e.g., google-safe-browsing-mcp, google-reputation-mcp).
- Other search engines (Bing, DuckDuckGo) deploy their own crawlers using Google's servers as dependencies.
- This doesn't mean Google shares raw data; it means standardized, audited APIs for safety/reputation checks.

4. Freshness optimization:

- MCP coordination server tracks "last seen" timestamp for each URL.
- Frontier server dynamically adjusts crawl frequency: popular sites crawled hourly, dormant sites crawled quarterly.
- LLM-based agents reason about freshness priorities (news sites need hourly checks, documentation sites need weekly).

Benefits:

- Cleaner architecture, easier to evolve and extend.
- Can incorporate external reputation signals (e.g., Perplexity's credibility scores) without reimplementing crawling logic.
- Search quality improves by incorporating richer semantic signals (not just links).

3.3 Microsoft: Copilot's Enterprise Crawling & Compliance

Current Challenge: Copilot needs to crawl enterprise intranets and knowledge bases while respecting security, authentication, and compliance (GDPR, HIPAA). Each enterprise has different policies.

MGWC Solution:

1. Deploy enterprise-specific policy servers:

- Each customer (e.g., a financial services firm) hosts an MCP Policy Server exposing:
 - check_crawl_compliance: "Is it OK to crawl this document? Is the user authorized?"
 - get_compliance_rules: "What data anonymization rules apply?"
- Microsoft's Copilot agents can crawl enterprise data without hardcoding enterprise logic.

2. Federated authentication:

- MCP auth layer supports OAuth, SAML, Kerberos, custom tokens.
- When Guide Bot crawls enterprise site, it authenticates via enterprise's OAuth provider.

- No credential sharing; Copilot agents only see what the authenticated user is allowed to see.

3. Crawl auditing & compliance:

- Every crawl decision logged via MCP audit tools.
- Compliance teams can query: "Which documents did the Copilot agent crawl and why?"
- Automatic retention policies (e.g., "Forget this document after 30 days per GDPR").

4. Integration with enterprise knowledge graphs:

- Crawled content automatically integrated with enterprise knowledge graphs (organizational structure, product catalogs, policies).
- Copilot agents reason over unified knowledge: "What's our policy on cloud pricing?" → crawls internal wiki, integrates with vendor contracts stored in Graph.

Benefits:

- Enterprise customers trust Copilot more (transparent, auditable).
- Microsoft can scale to thousands of enterprises without custom integrations.
- Compliance & legal teams have visibility and control.

3.4 Perplexity: Specialized, Real-Time, Multi-Source Crawling

Current Challenge: Perplexity's value proposition is answering questions with real-time, multi-source citations. But maintaining crawl freshness across diverse sources (news, academic papers, product reviews, Reddit discussions) is complex.

MGWC Solution:

1. Deploy specialized Guide Bots:

- Reddit crawler: Crawls Reddit discussions, forums, Q&A sites (updated every 4 hours).
- News crawler: Reuters, AP, BBC, local news sites (updated every 30 minutes).
- Academic crawler: ArXiv, Google Scholar, ResearchGate (updated daily).
- Product review crawler: Trustpilot, Amazon reviews, specialized review sites (updated hourly).

2. MCP frontier server prioritization:

- User asks: "What are the best running shoes for marathons?"
- Coordination server recognizes query needs product reviews + user discussions + expert recommendations.
- Triggers targeted crawls: Reddit (marathon running communities), review sites (Wirecutter, Runner's World), product sites (Nike, Brooks).

3. Real-time ranking & deduplication:

- MCP cache server detects if content already crawled by another bot or another query.
- Avoids crawling the same Reddit post twice in 10 minutes.
- Coordinates crawl rate with other Perplexity users' concurrent queries.

4. Trust & source scoring:

- MCP reputation server scores sources: Reddit comments get lower trust than expert reviews; established news agencies get highest trust.
- Answer generation prioritizes high-trust sources; lower-trust sources used for context/supplementation.

Benefits:

- More efficient crawling (less duplication, more targeted).
- Better answer quality (prioritizes high-trust sources).
- Faster user response time (crawl results can be cached and reused across similar queries).

3.5 Academic & Research Organizations: Distributed Knowledge Crawling

Current Challenge: Researchers at universities, nonprofits, and think tanks need to crawl large corpora (legal documents, scientific papers, government data) but lack the infrastructure of tech giants.

MGWC Solution:

1. Shared crawling infrastructure:

- Universities deploy lightweight MCP servers exposing their repositories (institutional knowledge bases, digital libraries).
- Other researchers' Guide Bots can crawl these repositories via MCP, with audit trails and compliance.

2. Distributed crawling network:

- Researchers at MIT, Stanford, Cambridge, and others deploy Guide Bots that crawl in parallel.
- MCP coordination server prevents duplication: Bot A crawls MIT's repository, Bot B crawls Stanford's, etc.
- Results shared back to all researchers via MCP resources.

3. Open science & collaboration:

- Researchers publish data and models as MCP servers ("MIT's Citation Network Server", "Stanford's Graph Analysis Tools").
- Other researchers build Guide Bots that consume these servers.

Benefits:

- Access to cutting-edge crawling techniques without reinventing.
- Reduced computational cost (shared infrastructure).
- Reproducibility (all crawls audited and logged via MCP).

4. MCP Server Specification for Web Crawling

This section details the **MCP tools, resources, and prompts** that must be implemented by compliant MGWC systems.

4.1 MCP Frontier Manager Server

Purpose: Manages the crawling frontier, URL prioritization, and seed generation.

Tools

Tool 1: `frontier.score_url`

Scores a URL for crawlability and relevance. Guide Bots query this before deciding to crawl.

```
{  
  "name": "frontierscore_url",  
  "description": "Score a candidate URL for crawlability, relevance, and priority. Returns a priority score (0-1) and reasoning.",  
  "inputSchema": {  
    "type": "object",  
    "properties": {  
      "url": {  
        "type": "string",  
        "description": "The URL to score (must be absolute, fully qualified)"  
      },  
      "source_page": {  
        "type": "string",  
        "description": "The URL from which this URL was discovered (for link context)"  
      },  
      "context": {  
        "type": "object",  
        "properties": {  
          "topic": {  
            "type": "string",  
            "description": "Crawl topic/focus (e.g., 'machine_learning', 'financial_news')"  
          },  
          "depth": {  
            "type": "integer",  
            "description": "Crawl depth from seed (0 = seed URL, 1 = links from seed, etc.)"  
          },  
          "last_crawl": {  
            "type": "string",  
            "format": "date-time",  
            "description": "When this URL was last crawled (if ever)"  
          },  
          "anchor_text": {  
            "type": "string",  
            "description": "The anchor text of the link (for semantic context)"  
          },  
          "page_title": {  
            "type": "string",  
            "description": "Title of the source page (if available)"  
          },  
          "domain_priority": {  
            "type": "string",  
            "enum": ["critical", "high", "medium", "low", "unknown"],  
            "description": "Pre-assigned priority of the domain"  
          }  
        }  
      },  
      "required": ["url"]  
    }  
  }  
}
```

Response:

```
{  
  "url": "https://arxiv.org/abs/2401.12345",  
  "score": 0.92,  
  "priority": "high",  
  "reasoning": "Recent AI research paper matching topic 'machine_learning'. Domain arxiv.org has critical priority. Last crawled 2 days ago; content likely updated.",  
  "tags": ["research_paper", "recent", "high_quality"],  
  "estimated_content_type": "academic_paper",  
  "freshness_urgency": "medium",  
  "suggested_retry_after": 604800  
}
```

Key fields:

- **score** (0-1): The frontier server's confidence this URL is worth crawling. Typically, crawlers crawl URLs with score > 0.5.
- **priority** (critical/high/medium/low): Guides crawl ordering. Multiple crawlers coordinate to crawl "critical" URLs first.
- **tags**: Semantic tags (e.g., "paywall", "login_required", "video", "broken_link"). Guide Bots use tags to decide handling.
- **suggested_retry_after**: Seconds before this URL should be crawled again. If the URL is infrequently updated, suggest long retry periods.

Rationale: This tool is the core of MGWC. Instead of a fixed heuristic (e.g., "crawl if depth < 5"), the frontier server can apply sophisticated logic: LLM-based relevance, domain reputation, freshness requirements, and user-defined policies.

Tool 2: frontier.get_next_seeds

Returns a batch of high-quality seed URLs for a crawl campaign.

```
{  
  "name": "frontierget_next_seeds",  
  "description": "Get the next batch of seed URLs to start a crawl. Used when a crawler's local frontier is exhausted.",  
  "inputSchema": {  
    "type": "object",  
    "properties": {  
      "topic": {  
        "type": "string",  
        "description": "Crawl topic (e.g., 'financial_news', 'academic_research')"  
      },  
      "num_seeds": {  
        "type": "integer",  
        "minimum": 1,  
        "maximum": 1000,  
        "description": "Number of seed URLs to return"  
      },  
      "crawl_depth": {  
        "type": "integer",  
        "description": "Depth of crawl to perform"  
      }  
    }  
  }  
}
```

```

"description": "Maximum depth the crawler is willing to crawl from seeds"
},
"exclusions": {
"type": "array",
"items": { "type": "string" },
"description": "Domains to exclude from seeds (e.g., ['internal.company.com'])"
},
"preferences": {
"type": "object",
"properties": {
"prefer_high_authority": { "type": "boolean" },
"prefer_recent": { "type": "boolean" },
"prefer_diverse_sources": { "type": "boolean" },
"content_types": {
"type": "array",
"items": { "type": "string" },
"description": "Preferred content types: ['html', 'pdf', 'json', 'xml', 'video']"
}
}
},
"required": ["topic", "num_seeds"]
}
}

```

Response:

```

{
"seeds": [
{
"url": "https://www.bloomberg.com/markets",
"seed_rank": 1,
"authority_score": 0.99,
"rationale": "Bloomberg Markets is the top authority for financial news."
},
{
"url": "https://www.cnbc.com/markets",
"seed_rank": 2,
"authority_score": 0.98,
"rationale": "CNBC Markets is a trusted financial news source."
},
{
"url": "https://finance.yahoo.com",
"seed_rank": 3,
"authority_score": 0.95,
"rationale": "Yahoo Finance aggregates stock data and news."
}
],
"total_available_seeds": 150,
"refresh_interval_seconds": 3600,
"notes": "These seeds are selected based on authority and recency for topic"

```

```
'financial_news'."  
}
```

Rationale: Rather than making Guide Bots find seeds themselves, the MCP frontier server provides curated, high-quality seeds. Reduces cold start and improves crawl efficiency.

Tool 3: frontier.report_seen

Guide Bots report URLs they've crawled, allowing the frontier server to track global coverage and deduplicate.

```
{  
  "name": "frontier:report_seen",  
  "description": "Report that a URL has been crawled. Used for deduplication and global tracking.",  
  "inputSchema": {  
    "type": "object",  
    "properties": {  
      "crawler_id": {  
        "type": "string",  
        "description": "Unique ID of the crawler making the report"  
      },  
      "urls_crawled": {  
        "type": "array",  
        "items": {  
          "type": "object",  
          "properties": {  
            "url": { "type": "string" },  
            "content_hash": {  
              "type": "string",  
              "description": "SHA-256 hash of page content (for detecting near-duplicates)"  
            },  
            "status_code": { "type": "integer" },  
            "crawl_time": { "type": "string", "format": "date-time" }  
          }  
        }  
      },  
      "required": ["crawler_id", "urls_crawled"]  
    }  
  }
```

Response:

```
{  
  "status": "acknowledged",  
  "urls_deduplicated": 3,  
  "urls_new": 47,  
  "global_coverage_status": {  
    "total_urls_known": 5000000,  
    "total_urls_crawled": 2500000,  
    "coverage_percentage": 50.0
```

```
}
```

Rationale: Enables global deduplication. If two crawlers both discovered the same URL, the frontier server tracks it once.

Tool 4: frontier.get_domain_policy

Retrieves crawl policies for a specific domain (rate limits, freshness requirements, allowed/disallowed paths).

```
{
  "name": "frontier.get_domain_policy",
  "description": "Get crawl policies for a domain (rate limits, allowed/disallowed paths, freshness).",
  "inputSchema": {
    "type": "object",
    "properties": {
      "domain": {
        "type": "string",
        "description": "Domain name (e.g., 'example.com')"
      }
    },
    "required": ["domain"]
  }
}
```

Response:

```
{
  "domain": "arxiv.org",
  "policies": {
    "rate_limit_requests_per_second": 1,
    "rate_limit_burst_size": 5,
    "freshness_requirements": {
      "recent_content_refresh_every_hours": 24,
      "stale_content_refresh_every_hours": 168
    },
    "disallowed_paths": ["/admin", "/api/internal"],
    "required_user_agent": "MCP-Crawler",
    "robots_txt_cache": {
      "last_updated": "2025-12-16T09:00:00Z",
      "ttl_seconds": 86400
    },
    "source": "robots.txt + domain_policy_database"
  }
}
```

Rationale: Ensures crawlers respect domain policies without re-parsing robots.txt for every request.

Resources

Resource 1: frontier/topic_taxonomy

A resource exposing the list of available crawl topics and their definitions.

```
{  
  "name": "frontier/topic_taxonomy",  
  "description": "Available crawl topics and their definitions",  
  "mimeType": "application/json",  
  "contents": {  
    "topics": {  
      "financial_news": {  
        "description": "Financial markets, earnings, economic indicators",  
        "priority_domains": ["bloomberg.com", "reuters.com", "cnbc.com"],  
        "excluded_domains": ["conspiracy_financial_sites.txt"]  
      },  
      "academic_research": {  
        "description": "Peer-reviewed papers, conferences, preprints",  
        "priority_domains": ["arxiv.org", "scholargoogle.com", "ieee.org"]  
      },  
      "e_commerce": {  
        "description": "Product listings, reviews, pricing",  
        "priority_domains": ["amazon.com", "ebay.com", "etsy.com"]  
      }  
    }  
  }  
}
```

Rationale: Guide Bots discover available topics without hardcoding them.

4.2 MCP Reputation Server

Purpose: Maintains and shares domain/URL reputation scores, trust signals, and safety information.

Tools

Tool 1: reputation.get_domain_score

Returns a domain's reputation score and trust factors.

```
{  
  "name": "reputation.get_domain_score",  
  "description": "Get reputation score for a domain, including trust factors and warnings.",  
  "inputSchema": {  
    "type": "object",  
    "properties": {  
      "domain": {  
        "type": "string",  
        "description": "Domain name"  
      },  
      "include_history": {  
        "type": "boolean",  
        "description": "Include historical data"  
      }  
    }  
  }  
}
```

```
"type": "boolean",
"description": "Include historical score trends"
}
},
"required": ["domain"]
}
}
```

Response:

```
{
"domain": "arxiv.org",
"trust_score": 0.99,
"trust_score_breakdown": {
"authority": 0.99,
"content_quality": 0.98,
"uptime_reliability": 0.99,
"security": 0.99,
"legal_compliance": 0.95
},
"signals": {
"is_https": true,
"has_valid_cert": true,
"whois_age_years": 25,
"organization": "Cornell University",
"page_rank_estimate": 8.5,
"spam_likelihood": 0.01,
"malware_likelihood": 0.0,
"phishing_likelihood": 0.0
},
"warnings": [],
"last_updated": "2025-12-16T08:00:00Z"
}
```

Rationale: Allows crawlers to make trust-based crawling decisions. High-trust domains might be crawled more aggressively; low-trust domains avoided entirely.

Tool 2: reputation.check_url_safety

Checks a specific URL for malware, phishing, or policy violations.

```
{
"name": "reputation.check_url_safety",
"description": "Check if a URL is safe to visit (no malware, phishing, policy violations).",
"inputSchema": {
"type": "object",
"properties": {
"url": {
"type": "string"
},
"check_types": {
"type": "array",

```

```

"items": { "type": "string" },
"enum": ["malware", "phishing", "nsfw", "copyright_violation", "policyViolation"],
"description": "Types of checks to perform"
}
},
"required": ["url"]
}
}

```

Response:

```

{
"url": "https://example.com/article",
"is_safe": true,
"checks": {
"malware": {
"status": "safe",
"confidence": 0.99,
"last_scanned": "2025-12-16T07:00:00Z"
},
"phishing": {
"status": "safe",
"confidence": 0.98
},
"nsfw": {
"status": "safe",
"confidence": 0.95,
"likely_content_type": "technology_blog"
}
},
"recommendations": "Safe to crawl. Content is appropriate for public indexing."
}

```

Rationale: Prevents crawlers from visiting malicious or inappropriate content. Can integrate with Google Safe Browsing, Cisco Talos, or other threat feeds.

Tool 3: reputation.report_malicious

Guide Bots report URLs that they've identified as malicious or policy-violating.

```

{
"name": "reputation.report_malicious",
"description": "Report a URL as malicious or policy-violating.",
"inputSchema": {
"type": "object",
"properties": {
"url": { "type": "string" },
"issue_type": {
"type": "string",
"enum": ["malware", "phishing", "copyrighted_content", "hate_speech", "illegal_activity"]
},
"confidence": {

```

```
"type": "number",
"minimum": 0,
"maximum": 1
},
"evidence": {
"type": "string",
"description": "Description of evidence or context"
}
},
"required": ["url", "issue_type", "confidence"]
}
}
```

Response:

```
{
"status": "received",
"url": "https://suspicious.site/malware.exe",
"issue_type": "malware",
"confidence": 0.92,
"alert_level": "high",
"action_taken": "URL added to blocklist. All crawlers will receive warning."
}
```

Rationale: Enables collective threat intelligence. One crawler's discovery is shared with all others.

4.3 MCP Coordination Server

Purpose: Orchestrates multi-crawler campaigns, prevents duplication, manages crawler registry and assignments.

Tools

Tool 1: coordination.register_crawler

A Guide Bot registers itself with the coordination server when it starts.

```
{
"name": "coordination.register_crawler",
"description": "Register a crawler instance with the coordination server",
"inputSchema": {
"type": "object",
"properties": {
"crawler_id": {
"type": "string",
"description": "Unique ID for this crawler"
},
"organization": {
"type": "string",
"description": "Organization running the crawler (e.g., 'OpenAI')"
}
},
```

```

"crawler_version": {
  "type": "string"
},
"topic": {
  "type": "string",
  "description": "Crawl topic (e.g., 'financial_news')"
},
"capabilities": {
  "type": "array",
  "items": { "type": "string" },
  "description": "Special capabilities (e.g., ['javascript_rendering', 'video_download', 'authentication'])"
},
"geographic_focus": {
  "type": "array",
  "items": { "type": "string" },
  "description": "Geographic regions this crawler focuses on (e.g., ['US', 'EU'])"
}
},
"required": ["crawler_id", "topic"]
}
}

```

Response:

```

{
  "status": "registered",
  "crawler_id": "guidebot_001",
  "assigned_work": {
    "domain_assignments": ["techcrunch.com", "wired.com"],
    "topic": "tech_news",
    "priority_level": "high"
  },
  "coordination_server": "coord.example.com:9000",
  "mcp_servers_available": [
    "frontier.mcp.example.com",
    "reputation.mcp.example.com",
    "cache.mcp.example.com"
  ]
}

```

Rationale: Allows coordination server to track active crawlers, balance load, and assign work.

Tool 2: coordination.report_crawl_status

Crawlers periodically report progress.

```
{
  "name": "coordination.report_crawl_status",
  "description": "Report progress on the crawl job.",
  "inputSchema": {
    "type": "object",
    "properties": {
      "crawl_id": { "type": "string" },
      "progress": { "type": "number" }
    }
  }
}
```

```

"type": "object",
"properties": {
  "crawler_id": { "type": "string" },
  "url_crawled": { "type": "string" },
  "status": {
    "type": "string",
    "enum": ["success", "error", "rate_limited", "blocked", "invalid_url"]
  },
  "num_new_urls_discovered": { "type": "integer" },
  "content_hash": { "type": "string" },
  "error_message": { "type": "string" },
  "response_time_ms": { "type": "integer" }
},
"required": ["crawler_id", "url_crawled", "status"]
}
}

```

Response:

```
{
  "status": "acknowledged",
  "crawler_message": "Keep going. You're doing great.",
  "global_coverage_status": {
    "urls_crawled_total": 2500000,
    "urls_pending": 500000,
    "estimated_completion_hours": 24
  },
  "rate_limit_adjustment": "No change"
}
```

Rationale: Provides real-time visibility and allows the coordination server to issue dynamic guidance (e.g., slow down, speed up, shift focus).

Tool 3: coordination.sync_seen_urls

Crawlers share the URLs they've seen, enabling global deduplication.

```
{
  "name": "coordination.sync_seen_urls",
  "description": "Sync the set of URLs this crawler has seen with the global registry.",
  "inputSchema": {
    "type": "object",
    "properties": {
      "crawler_id": { "type": "string" },
      "urls_seen": {
        "type": "array",
        "items": { "type": "string" },
        "description": "List of URLs crawled or skipped by this crawler"
      },
      "url_bloom_filter": {
        "type": "string",
        "description": "Compressed Bloom filter of all URLs seen (for large-scale syncs)"
      }
    }
  }
}
```

```
    },
    },
    "required": ["crawler_id"]
}
}
```

Response:

```
{
  "status": "synced",
  "new_urls_from_other_crawlers": [
    "https://example.com/new-article",
    "https://another-site.com/page"
  ],
  "global_seen_count": 5000000,
  "dedup_efficiency": 0.92
}
```

Rationale: Enables cross-crawler deduplication. If Crawler A just learned about a URL, Crawler B can fetch it from this response rather than discovering it independently.

4.4 MCP Policy Enforcement Server

Purpose: Enforces crawl policies, compliance rules, and respect for domain boundaries.

Tools

Tool 1: `policy.check_crawl_compliance`

Check if crawling a URL complies with all applicable policies.

```
{
  "name": "policy.check_crawl_compliance",
  "description": "Check if crawling a URL is compliant with all policies (robots.txt, GDPR, organization policy, etc.)",
  "inputSchema": {
    "type": "object",
    "properties": {
      "url": { "type": "string" },
      "domain": { "type": "string" },
      "context": {
        "type": "object",
        "properties": {
          "user_id": { "type": "string" },
          "organization": { "type": "string" },
          "crawl_purpose": { "type": "string" },
          "jurisdictions": { "type": "array", "items": { "type": "string" } }
        }
      }
    },
    "required": ["url"]
```

```
}
```

```
}
```

Response:

```
{
  "url": "https://example.com/page",
  "is_compliant": true,
  "checks": {
    "robots_txt": {
      "status": "allowed",
      "rule": "User-agent: MCP-Crawler\\nAllow: /"
    },
    "gdpr": {
      "status": "compliant",
      "notes": "URL is public, no PII detected"
    },
    "organization_policy": {
      "status": "allowed",
      "reason": "Example.com is on whitelist"
    }
  },
  "warnings": []
}
```

Rationale: Prevents legal/compliance violations. Especially important for enterprise crawling where GDPR, CCPA, and organizational policies apply.

Tool 2: policy.get_rate_limits

Retrieve rate limits for a domain, respecting both domain-side policies (robots.txt) and organizational policies.

```
{
  "name": "policy.get_rate_limits",
  "description": "Get the allowed rate limits for crawling a domain.",
  "inputSchema": {
    "type": "object",
    "properties": {
      "domain": { "type": "string" },
      "organization": { "type": "string" }
    },
    "required": ["domain"]
  }
}
```

Response:

```
{
  "domain": "example.com",
  "rate_limits": {
    "requests_per_second": 0.5,
```

```
"burst_size": 3,  
"concurrent_connections": 2  
},  
"source": ["robots.txt", "domain_policy_db"],  
"coordination_notes": "This domain is also being crawled by 2 other crawlers. Globally, rate  
is set to 0.167 req/s per crawler to total 0.5 req/s."  
}
```

Rationale: Ensures crawlers respect domain limits and coordinate globally to avoid overwhelming servers.

4.5 MCP Safety Gateway Server

Purpose: Scans URLs and content for malware, phishing, NSFW content, and other unsafe material.

Tools

Tool 1: safety.scan_url_for_threats

Scan a URL for malware, phishing, and other threats.

```
{  
"name": "safety.scan_url_for_threats",  
"description": "Scan a URL for malware, phishing, ransomware, and other threats.",  
"inputSchema": {  
"type": "object",  
"properties": {  
"url": { "type": "string" },  
"use_advanced_heuristics": { "type": "boolean" }  
},  
"required": ["url"]  
}  
}
```

Response:

```
{  
"url": "https://example.com",  
"threat_level": "safe",  
"threats": {  
"malware": { "status": "not_detected", "confidence": 0.99 },  
"phishing": { "status": "not_detected", "confidence": 0.98 },  
"ransomware": { "status": "not_detected", "confidence": 0.99 }  
},  
"recommendation": "Safe to crawl"  
}
```

Tool 2: safety.check_content_policy

Check if content violates organizational or legal content policies (NSFW, hate speech, violence, etc.).

```
{
  "name": "safety.check_content_policy",
  "description": "Check if content violates content policies (NSFW, hate speech, violence, etc.)",
  "inputSchema": {
    "type": "object",
    "properties": {
      "content": {
        "type": "string",
        "description": "Content (text or HTML) to check"
      },
      "policy_categories": {
        "type": "array",
        "items": { "type": "string" },
        "enum": ["adult", "violence", "hate_speech", "misinformation", "copyright"]
      }
    },
    "required": ["content"]
  }
}
```

Response:

```
{
  "content_hash": "abc123",
  "policy_violations": {
    "adult": { "detected": false, "confidence": 0.99 },
    "violence": { "detected": false, "confidence": 0.95 },
    "hate_speech": { "detected": false, "confidence": 0.97 }
  },
  "overall_safe": true,
  "recommendation": "Content is suitable for indexing"
}
```

4.6 MCP Cache Server

Purpose: Caches frequently needed data (robots.txt, content signatures, URL fingerprints) to reduce redundant checks.

Tools

Tool 1: cache.lookup_duplicate

Check if a URL (or a near-duplicate) has been crawled before.

```
{
  "name": "cache.lookup_duplicate",
  "description": "Look up if a URL or near-duplicate has been seen before.",
  "inputSchema": {
    "type": "object",
    "properties": {
      "url": { "type": "string" },
      "threshold": { "type": "number", "minimum": 0, "maximum": 1 }
    }
  }
}
```

```
"content_signature": {  
    "type": "string",  
    "description": "Optional: pre-computed signature of content"  
},  
"check_near_duplicates": { "type": "boolean" }  
},  
"required": ["url"]  
}  
}
```

Response:

```
{  
    "url": "https://example.com/article",  
    "exact_match_found": true,  
    "previously_crawled": {  
        "timestamp": "2025-12-15T10:00:00Z",  
        "crawler_id": "guidebot_003",  
        "content_hash": "abc123xyz"  
    },  
    "recommendation": "Skip this URL (exact duplicate)"  
}
```

Tool 2: cache.store_content_signature

Store the signature of a crawled page for future duplicate detection.

```
{  
    "name": "cache.store_content_signature",  
    "description": "Store the content signature of a crawled page.",  
    "inputSchema": {  
        "type": "object",  
        "properties": {  
            "url": { "type": "string" },  
            "content_hash": { "type": "string" },  
            "crawler_id": { "type": "string" },  
            "crawl_time": { "type": "string", "format": "date-time" }  
        },  
        "required": ["url", "content_hash", "crawler_id"]  
    }  
}
```

4.7 Authentication & Authorization

All MCP servers must support authentication. Recommended approach:

MCP Authentication Profile for Crawling:

- Transport: TLS 1.3+
- Authentication Methods:
 - OAuth 2.0 (for organization-to-organization federation)
 - Mutual TLS (for trusted service-to-service)

- API Keys with rotation (for development)
 - SAML 2.0 (for enterprise environments)
 - Authorization Model (RBAC):
 - Role: crawler (can call frontier, reputation, coordination tools)
 - Role: coordinator (can call all tools, plus policy overrides)
 - Role: admin (full access, can modify policies)
 - Audit Logging:
 - Every MCP call logged with: timestamp, caller_id, tool_name, input_hash, output_hash, decision (allow/deny)
 - Audit logs queryable via MCP resource: audit/recent_calls
-

5. Agent Coordination Patterns

5.1 The Guide Bot Agent Architecture

A Guide Bot is an **autonomous agent** (powered by an LLM or agent framework) that:

1. **Reads** crawl objectives from a Campaign Manager.
2. **Plans** a crawl strategy (which seeds, which topics, which crawl depth).
3. **Executes** by fetching URLs, extracting links, and querying MCP servers.
4. **Adapts** in real-time based on MCP feedback (adjust scope, change strategies).
5. **Reports** progress and discoveries to coordination servers.

Pseudo-code of a Guide Bot's main loop:

```
class GuideBot(Agent):
    def __init__(self, mcp_client, campaign_objective):
        self.mcp = mcp_client
        self.objective = campaign_objective # "Index all academic AI research, 2023-2025"
        self.local_frontier = []
        self.crawl_stats = {}

    def run(self):
        # Step 1: Register with coordination server
        self.register()

        # Step 2: Get initial seeds
        seeds = self.mcp.call_tool("frontier.get_next_seeds", {
            "topic": self.objective.topic,
            "num_seeds": 10
        })
        self.local_frontier.extend(seeds)

        # Step 3: Main crawl loop
        while not self.should_stop():
            url = self.local_frontier.pop(0)
            # Process URL and update stats
            self.crawl_stats[url] = {"status": "processing", "progress": 0}
            # Add discovered URLs back to frontier
            self.local_frontier.extend(discovered_urls)
```

```
# Check if we've already crawled this
duplicate_check = self.mcp.call_tool("cache.lookup_duplicate", {"url": url})
if duplicate_check.get("exact_match_found"):
    continue # Skip duplicates

# Check compliance
compliance = self.mcp.call_tool("policy.check_crawl_compliance", {"url": url})
if not compliance["is_compliant"]:
    self.log(f"Skipping {url}: non-compliant")
    continue

# Check safety
safety = self.mcp.call_tool("safety.scan_url_for_threats", {"url": url})
if safety["threat_level"] == "dangerous":
    self.log(f"Skipping {url}: threat detected")
    continue

# Fetch and parse
try:
    html = self.fetch(url)
    links = self.extract_links(html)
except Exception as e:
    self.report_error(url, e)
    continue

# Score new links
for link in links:
    score_resp = self.mcp.call_tool("frontier.score_url", {
        "url": link,
        "source_page": url,
        "context": {"topic": self.objective.topic}
    })

    if score_resp["score"] > 0.6: # Threshold
        self.local_frontier.append(link)

# Report progress
```

```

        self.mcp.call_tool("coordination.report_crawl_status", {
            "crawler_id": self.id,
            "url_crawled": url,
            "status": "success",
            "num_new_urls_discovered": len(links)
        })

        # Periodically sync seen URLs
        if len(self.crawl_stats) % 100 == 0:
            self.mcp.call_tool("coordination.sync_seen_urls", {
                "crawler_id": self.id,
                "urls_seen": list(self.crawl_stats.keys())
            })

    def should_stop(self):
        # Stop if frontier is empty and no new seeds available
        if self.local_frontier:
            return False

        new_seeds = self.mcp.call_tool("frontier.get_next_seeds", {
            "topic": self.objective.topic,
            "num_seeds": 10
        })

        if new_seeds:
            self.local_frontier.extend(new_seeds)
            return False

    return True # No more work

```

5.2 Multi-Agent Collaboration

When multiple Guide Bots crawl concurrently, collaboration is essential:

Deduplication via Coordination Server:

Bot A: "I'm going to crawl <https://example.com>"

Coordination Server: "Roger. I'll mark it as claimed by Bot A. Bot B and C, don't crawl this."

Bot B: "Can I crawl <https://another-site.com>?"

Coordination Server: "Sure, that's available."

Bot A finishes <https://example.com>, reports 50 new URLs

Coordination Server: Adds these to global frontier, makes available to all bots

Bot C: "Got any work for me?"

Coordination Server: "Yes, here are 50 URLs from Bot A's recent crawl."

Load Balancing:

Initially, Bot A, B, C all crawling finance domain at 0.5 req/s each = 1.5 req/s total (exceeds policy)

Coordination Server detects this, instructs each bot to crawl at 0.167 req/s = 0.5 req/s total (compliant)

Failure Recovery:

Bot A crashes while crawling example.com

Coordination Server: "Bot A hasn't reported in 5 minutes. Reassigning its URLs to Bot B and C."

Bot B/C: Resume crawling from the frontier, no duplicate work lost.

6. Implementation Roadmap

Phase 1: Single-Crawler Integration (Weeks 1-4)

Goal: Prove the architecture works with one Guide Bot querying MCP servers.

Deliverables:

1. Deploy MCP Frontier Server (prototype)

- Simple REST API (or native MCP) exposing frontierscore_url and frontierget_next_seeds.
- Backed by a curated list of domains and topics.
- No ML-based scoring initially; use heuristics (domain authority, recency).

2. Adapt existing crawler to use MCP

- Modify Apache Nutch, Scrapy, or internal crawler to call frontierscore_url before crawling.
- Add configuration to specify MCP server endpoints.
- Measure: URLs crawled, average crawl quality, frontier latency.

3. MCP Reputation Server (simple prototype)

- Expose reputation.get_domain_score backed by pre-computed scores (e.g., from Majestic, Moz, or internal rankings).
- Initially, read-only (no threat intelligence integration yet).

4. Monitoring & Metrics

- Track: crawl speed, cache hit rate, MCP server latency, URLs per unique host.
 - Set baselines for comparison in later phases.
-

Phase 2: Multi-Agent Coordination (Weeks 5-12)

Goal: Deploy 3–5 Guide Bots coordinating via MCP coordination server.

Deliverables:

1. MCP Coordination Server

- Implement coordination.register_crawler, coordination.report_crawl_status, coordination.sync_seen_urls.
- Central database (Redis or PostgreSQL) tracking global state.
- API for querying coverage, active crawlers, pending work.

2. Distributed Deduplication

- Implement global URL deduplication (Bloom filter or hash set).
- Guide Bots report seen URLs; coordination server merges and broadcasts back.
- Measure: reduction in duplicate fetches, crawl efficiency improvement.

3. Multi-Topic Crawling

- Deploy specialized Guide Bots: finance_bot, tech_bot, research_bot.
- Each bot queries appropriate MCP servers (finance frontier server vs. research frontier server).
- Measure: topic-specific coverage, crawl relevance.

4. Load Balancing

- Implement rate limit coordination: coordination server instructs bots on per-domain request rates.
- Measure: compliance with domain rate limits, no "thundering herd" issues.

Phase 3: Enterprise Integration (Weeks 13-20)

Goal: Integrate with production systems (search index, training pipelines, knowledge graphs).

Deliverables:

1. Policy Enforcement Server

- Implement policy.check_crawl_compliance integrating robots.txt, GDPR rules, organizational policies.
- Fine-tune compliance checks with legal and security teams.

2. Safety Gateway Server

- Integrate with Google Safe Browsing, malware scanners.
- Implement safety.scan_url_for_threats and safety.check_content_policy.
- Test false positive rates; tune thresholds.

3. Integration with Search Index

- Pipeline: Crawled content → parsed, deduplicated, indexed → searchable.
- Set up indexing jobs triggered by crawl completion.

4. Analytics & Observability

- Real-time dashboard: crawl progress, topic coverage, error rates, MCP latency.
- Structured logging for compliance audits.

Phase 4: Ecosystem & Standardization (Weeks 21+)

Goal: Publish MGWC as an open standard; enable cross-organization collaboration.

Deliverables:

1. Publish MCP Server Specification

- Release detailed specs for all MCP servers (frontier, reputation, coordination, policy, safety).
- Open-source reference implementations (Python, JavaScript, Go).

2. Interoperability Testing

- Ensure crawlers from different organizations can call each other's MCP servers.
- Federated crawling test: OpenAI bot + Google bot + Perplexity bot crawling together.

3. Community Contributions

- Invite other organizations to contribute MCP servers (e.g., "Microsoft Threat Intel Server", "AWS Compliance Server").
- Build ecosystem of specialized servers.

4. Performance Optimization

- Profile and optimize MCP call latency (target: <50ms for frontierscore_url).
 - Implement caching layers, batch processing for bulk operations.
-

7. Production Considerations

7.1 Scalability

Horizontal Scaling:

- **MCP Servers:** Stateless, can run in replicas behind a load balancer. Scale individual servers independently (e.g., 10 frontier replicas, 5 reputation replicas).
- **Crawlers:** Lightweight; scale from 1 to 1000s easily. Each crawler is a small process.
- **Coordination Server:** Central, but designed to be a thin dispatch layer. Use distributed cache (Redis Cluster) for state.

Vertical Scaling:

- Increase MCP server hardware (CPU, memory) for specialized processing (e.g., reputation scoring with ML inference).

Estimated Capacity (single organization, Google-scale):

- 1000 Guide Bots crawling concurrently.
- 10 billion URLs in frontier.
- 1000 requests/sec to MCP servers.
- Requires: 20-30 MCP server replicas, 500GB Redis cluster, 100TB PostgreSQL for global state.

7.2 Fault Tolerance & Recovery

MCP Server Failures:

- Crawlers have fallback logic: if coordination server is down, continue crawling with local frontier (advisory mode).
- Use health checks: crawlers detect dead MCP servers and retry with backoff.
- Geographic redundancy: deploy MCP servers in multiple data centers.

Crawler Failures:

- Coordination server detects inactive crawlers (no heartbeat in 5 minutes).
- Reassigns its URLs to healthy crawlers.
- Minimal data loss (only in-flight URLs; already-crawled URLs are safe).

Data Loss Prevention:

- All crawl results written to durable store immediately (not in memory).
- Global state (seen URLs, domain policies) replicated across regions.

7.3 Security & Authorization

Threats & Mitigations:

1. Malicious crawler misusing MCP:

- MCP auth: all crawlers must authenticate (API key, OAuth).
- Rate limiting per crawler: prevent one bot from monopolizing resources.
- Audit logging: track all MCP calls for forensics.

2. Rogue MCP server:

- Mutual TLS: MCP clients verify server certificates.
- Signed responses: MCP servers sign outputs (e.g., domain scores) so clients can verify authenticity.

3. Data exfiltration:

- Crawled content should not be stored on untrusted systems.
- Encrypted transport (TLS); encrypted storage (AES-256).
- Access controls: only authorized pipelines can read crawled content.

Compliance:

- **GDPR:** Crawled personal data is handled with care (minimization, retention limits, user rights).
- **CCPA:** Respect data subject requests (right to deletion, etc.).
- **Robots.txt:** Always honored.
- **Copyright:** Respect copyright; don't store copyrighted content (only metadata/snippets).

7.4 Observability & Monitoring

Metrics to Track:

- Crawl speed: URLs crawled per second (target: >1000 URLs/sec per organization)
- Crawl quality: % of crawled URLs that are high-quality (target: >80%)
- Coverage: % of target domains covered (target: >95%)
- Error rate: % of URLs that resulted in errors (target: <5%)
- Duplicate rate: % of URLs that were duplicates (target: <10%)
- MCP latency: p50, p95, p99 latency for MCP calls (target: <50ms p95)
- Compliance: % of crawls that violate robots.txt or policies (target: 0%)
- Safety: # of unsafe URLs visited (target: 0)

Logging:

- Structured logs: crawler_id, url, timestamp, action (crawled/skipped/error), reason
- MCP audit logs: mcp_server, tool_name, caller, timestamp, decision
- Compliance logs: crawls that respect/violate policies
- Security logs: authentication failures, rate limit violations

8. Specific Use Cases for Major Organizations

8.1 OpenAI: Real-Time Knowledge for GPT & ChatGPT

Objective: Enable ChatGPT to answer real-time questions ("What's the latest OpenAI news?") with current sources.

Implementation:

1. Deploy news + tech topic crawlers:

- Tech news crawler: TechCrunch, Hacker News, ArXiv, GitHub.
- General news crawler: Reuters, AP, BBC.
- Financial crawler: Bloomberg, S&P, Fed Reserve.

2. High-speed frontier server:

- Update feed every 30 minutes to 2 hours for breaking news.
- Frontier server prioritizes freshness (crawl recent links first).

3. Tight integration with knowledge retrieval:

- User query → agent retrieves from crawl cache + knowledge base.
- If cache is stale (>4 hours), trigger a targeted recrawl.

4. Citation & attribution:

- Every answer includes source URL and crawl timestamp.
- Users can verify by clicking through to original source.

Metrics:

- Recency: 99% of answers based on content <24 hours old.
- Accuracy: <5% hallucination rate on factual queries.
- User satisfaction: >90% of users prefer cited, real-time answers.

8.2 Google: Modernizing Googlebot with MCP

Objective: Evolve Googlebot to be more intelligent, ranking-aware, and compliant.

Implementation:

1. Decompose Googlebot into Guide Bots + MCP:

- Replace monolithic frontier scheduler with MCP frontier server.
- Each Googlebot query: frontier server (what to crawl), reputation server (how to rank), policy server (compliance).

2. Topical specialization:

- News crawler: crawl news sites hourly.
- E-commerce crawler: crawl product pages daily.
- Academic crawler: crawl journals weekly.

3. Better E-E-A-T integration:

- Reputation server incorporates Google's E-E-A-T signals.
- Frontier server prioritizes expert, authoritative sites.

4. Faster policy enforcement:

- robots.txt cached centrally (MCP cache server).
- Policy enforcement (GDPR, CCPA, copyright) coordinated globally.

5. Collaboration with other search engines:

- Publish reputation server: other search engines can use Google's safety/trust data.

- o Subscribe to other organizations' threat feeds.

Metrics:

- Search quality: CTR, dwell time on search results improves.
 - Crawl efficiency: bandwidth savings (less duplicate crawling, better prioritization).
 - Compliance: 100% robots.txt compliance, 0% GDPR violations.
-

8.3 Microsoft: Copilot Enterprise Crawling

Objective: Enable Copilot to crawl enterprise intranets safely, compliantly, and with proper authorization.

Implementation:

1. Deploy federated policy servers:

- o Each enterprise customer hosts an MCP policy server
- o Copilot crawlers query customer's policy server before accessing enterprise data.

2. Authentication & authorization:

- o Crawlers authenticate as specific users (e.g., "Copilot crawling on behalf of user john@company.com").
- o Policy server checks: "Is john authorized to access this document?"

3. Compliance & auditing:

- o Every crawl logged: which document, which user, which Copilot, timestamp.
- o Enterprises can audit: "What data did Microsoft access?" Answer: full audit trail.

4. Data minimization:

- o Crawlers only fetch documents user is authorized to see.
- o Snippets, summaries stored; raw documents not kept.

5. Retention policies:

- o Customer can specify: "Forget this document after 30 days (GDPR compliance)."
- o Cache server automatically purges after TTL.

Metrics:

- Enterprise adoption: # of enterprise customers using Copilot crawling.
 - Compliance score: % of crawls compliant with enterprise + regulatory policies.
 - User trust: % of users comfortable with Copilot accessing their enterprise data.
-

8.4 Perplexity: Fast, Cited, Multi-Source Answers

Objective: Deliver answers faster by caching and de-duplicating crawls across concurrent queries.

Implementation:

1. Query-triggered crawling:

- o User query: "Best running shoes for marathons"
- o Coordination server recognizes query intent: needs product reviews + user discussions.

- Triggers parallel crawls: Amazon, Wirecutter, Reddit, running forums.

2. Smart caching:

- If answer to similar query was generated in last 2 hours, reuse cached crawl results.
- Cache hit rate: 70% (reuse crawls across users).

3. Deduplication:

- Multiple users asking similar questions in parallel: crawl once, serve to all.
- Coordination server groups similar queries.

4. Trust-based ranking:

- Reputation server scores sources: Wirecutter (high) > Reddit comments (medium) > random blog (low).
- Answer generation prioritizes high-trust sources, uses low-trust for supplementation.

5. Real-time updates for trending queries:

- Popular query detected (e.g., "iPhone 16 release date").
- Trigger immediate recrawl to get freshest data.

Metrics:

- Time to answer: <500ms for cached queries, <2s for fresh crawls.
 - Source diversity: average answer cites 3-5 different high-quality sources.
 - User satisfaction: >85% users rate answers as "accurate and well-sourced."
-

9. Risk Mitigation & Compliance

9.1 Legal & Ethical Risks

Copyright:

- Do not crawl and store full copyrighted text.
- Store only: URLs, metadata, snippets (fair use).
- Respect Copyright headers (X-Robots-Tag: noindex).
- License agreement: if site requires permission, ask via DMCA/copyright contact.

GDPR & CCPA:

- Personal data: crawlers should avoid crawling sites with sensitive PII.
- User rights: implement "right to be forgotten" (delete crawled data on request).
- Data minimization: only crawl data necessary for the use case.
- Consent: if required by law, obtain consent before crawling.

Robots.txt & ToS:

- Always honor robots.txt.
- Respect website terms of service (read carefully before crawling).
- If site explicitly forbids crawling, don't crawl.

Government & Compliance:

- Some jurisdictions restrict crawling (e.g., some countries restrict export of technical data).
- Work with legal team to ensure compliance.

9.2 Technical Risks

Bot Detection & Blocking:

- Sites increasingly use anti-bot measures (Cloudflare, DataDome).
- Mitigation: rotate User-Agent strings, add delays between requests, use proxies.
- But be ethical: if site forbids crawling, respect it rather than circumvent.

Malware & Phishing:

- Crawled content may contain drive-by exploits, malware links.
- Mitigation: run content through sandboxed analysis, don't render untrusted JavaScript in production systems.

DDoS Abuse:

- Misconfigured crawlers can become DDoS vectors (e.g., crawl infinite loops).
- Mitigation: implement depth limits, size limits, timeout limits. Monitor crawl patterns for anomalies.

9.3 Operational Risks

Reputation Damage:

- Aggressive crawling can damage organizational reputation ("Amazonbot crawled our site 1 million times in an hour").
- Mitigation: test crawlers on staging domains first. Use polite rate limits. Monitor for complaints.

Dependency on MCP Servers:

- If coordination server goes down, crawlers can't dedup or coordinate.
- Mitigation: design crawlers to work in degraded mode (no coordination, but still functional). Use multiple coordination servers (redundancy).

10. Conclusion & Future Directions

10.1 Summary

MCP-Guided Web Crawling (MGWC) represents a fundamental shift in how organizations approach web indexing and data collection. By building crawling systems on the **Model Context Protocol**, we enable:

- **Agentic autonomy:** AI agents make intelligent, reasoned crawling decisions.
- **Standardization:** Multiple organizations interoperate seamlessly on shared infrastructure.
- **Scalability:** Distributed, fault-tolerant crawling from millions of concurrent agents.
- **Compliance & Safety:** Policies enforced centrally; trust and safety built-in.
- **Efficiency:** Global deduplication, smart prioritization, minimal redundant work.

The framework has been detailed with:

1. **Architectural foundations** explaining why MCP is the right abstraction.

2. **Real-world use cases** for OpenAI, Google, Microsoft, Perplexity, and academics.
3. **Production-ready MCP server specs** that implementers can follow immediately.
4. **Multi-agent coordination patterns** enabling large-scale, collaborative crawling.
5. **Implementation roadmap** from prototyping to standardization.
6. **Production hardening guidelines** for scalability, security, compliance.

10.2 Immediate Next Steps for Implementers

For OpenAI, Google, Microsoft, Perplexity, and other organizations:

1. **Form an MGWC working group** within your organization and externally with peers.
2. **Build a prototype:** Phase 1 crawler (single Guide Bot, basic MCP frontier server).
Timeline: 2-4 weeks.
3. **Iterate on real crawl workloads:** Does MGWC improve crawl quality and efficiency vs. current approach? Measure.
4. **Standardize on MCP:** Adopt the MCP server specs outlined in Section 4. Contribute feedback.
5. **Publish & collaborate:** Share your MCP servers (e.g., "Google Safe Browsing MCP") with the ecosystem. Enable cross-organizational crawling.

10.3 Research Opportunities

- **ML-based frontier scoring:** Use LLMs to predict page relevance, quality, and freshness more accurately.
- **Federated learning for reputation:** Train reputation models across organizations without sharing raw data.
- **Decentralized coordination:** Blockchain or DHT-based coordination for truly peer-to-peer crawling (no central server).
- **Energy efficiency:** Optimize crawl strategies to minimize energy consumption (important at global scale).

10.4 Conclusion

The future of web crawling is **agentic, distributed, and standardized**. The Model Context Protocol provides the foundation. It's time for the industry—OpenAI, Google, Microsoft, Perplexity, and others—to build the next generation of web crawlers together.

Appendices

Appendix A: MCP Protocol Primer

For readers unfamiliar with MCP, here's a quick overview:

What is MCP?

MCP is an open standard (JSON-RPC 2.0) for AI agents to call functions (tools) exposed by external services (servers). Think of it as "APIs for AI agents."

Basic Structure:

```
// Client calls tool on server
{
  "jsonrpc": "2.0",
  "method": "get_index",
  "params": {
    "url": "https://www.example.com"
  },
  "id": 1
}
```

```

"id": 1,
"method": "tools/call",
"params": {
"name": "frontierscore_url",
"arguments": {
"url": "https://example.com",
"topic": "tech_news"
}
}
}

// Server responds
{
"jsonrpc": "2.0",
"id": 1,
"result": {
"score": 0.92,
"priority": "high",
"reasoning": "..."
}
}

```

Key Concepts:

- **Tools:** Functions that MCP servers expose (e.g., frontierscore_url).
- **Resources:** Data that servers provide (e.g., frontier/topic_taxonomy).
- **Prompts:** System messages that guide agent behavior (e.g., "You are a web crawler assistant").

For detailed MCP documentation, see: <https://modelcontextprotocol.io/>

Appendix B: Sample Configuration File

Here's an example MGWC configuration for a Guide Bot:

config.yaml for Guide Bot instance

```

guidebot:
id: "guidebot_finance_001"
organization: "OpenAI"
version: "1.0.0"
topic: "financial_news"

crawl_parameters:
max_depth: 5
max_urls_per_domain: 10000
request_timeout_seconds: 10
max_concurrent_requests: 5

mcp_servers:
frontier:

```

```
url: "frontiermcp.example.com:9000"
auth_method: "oauth2"
auth_token_file: "/etc/mgwc/frontier_token"

reputation:
url: "reputation.mcp.example.com:9001"
auth_method: "mutual_tls"
cert_file: "/etc/mgwc/reputation_client.crt"
key_file: "/etc/mgwc/reputation_client.key"

coordination:
url: "coordination.mcp.example.com:9002"
auth_method: "api_key"
api_key_file: "/etc/mgwc/coordination_key"

policy:
url: "policy.mcp.example.com:9003"
auth_method: "oauth2"
auth_token_file: "/etc/mgwc/policy_token"

safety:
url: "safety.mcp.example.com:9004"
auth_method: "oauth2"
auth_token_file: "/etc/mgwc/safety_token"

crawl_policies:
respect_robots_txt: true
max_requests_per_second_per_domain: 0.5
crawl_delay_seconds: 1
user_agent: "MCP-GuideBot/1.0 (+https://mgwc.example.com)"

output:
index_endpoint: "https://search-index.example.com/api/index"
index_batch_size: 100
index_batch_timeout_seconds: 30

monitoring:
prometheus_endpoint: "localhost:9090"
log_level: "info"
log_file: "/var/log/guidebot.log"
```

Appendix C: MCP Server Implementation Example (Python)

Here's a minimal example of implementing a Frontier MCP server in Python:

frontier_server.py

Minimal MCP frontier server using Anthropic's Python SDK

```
from mcp.server import Server
from mcp.server:stdio import stdio_server
import json

app = Server("mcp-frontier-server")
```

Simulated domain scores and seed database

```
DOMAIN_SCORES = {
    "arxiv.org": 0.99,
    "techcrunch.com": 0.95,
    "medium.com": 0.75,
    "spam-site.com": 0.1,
}

SEEDS_BY_TOPIC = {
    "academic_research": [
        "https://arxiv.org/list/cs.AI/recent",
        "https://scholargoogle.com/scholar?q=artificial+intelligence",
    ],
    "tech_news": [
        "https://techcrunch.com/",
        "https://news.ycombinator.com/",
    ],
}

@app.call_tool()
async def score_url(url: str, source_page: str = None, context: dict = None):
    """Score a URL for crawlability."""

    # Extract domain from URL
    from urllib.parse import urlparse
    domain = urlparse(url).netloc

    # Lookup domain score
    base_score = DOMAIN_SCORES.get(domain, 0.5)

    # Adjust based on context
    if context:
        topic = context.get("topic", "")
```

```

# Bonus for on-topic links
if "academic" in topic and "arxiv" in domain:
    base_score = min(base_score + 0.1, 1.0)

return {
    "url": url,
    "score": base_score,
    "priority": "high" if base_score > 0.8 else "medium" if base_score > 0.5 else "low",
    "reasoning": f"Domain {domain} has base score {base_score}",
}

```

```

@app.call_tool()
async def get_next_seeds(topic: str, num_seeds: int = 10, **kwargs):
    """Get seed URLs for a crawl topic."""

```

```

seeds = SEEDS_BY_TOPIC.get(topic, [])

return {
    "seeds": [
        {"url": s, "seed_rank": i+1, "authority_score": 0.9}
        for i, s in enumerate(seeds[:num_seeds])
    ],
    "total_available_seeds": len(seeds),
}

```

```

if name == "main":
    stdio_server(app).run()

```

To run this server:

```

pip install mcp
python frontier_server.py

```

Clients (Guide Bots) can now call this server via MCP.

Appendix D: Further Reading & References

1. **Model Context Protocol Specification:** <https://modelcontextprotocol.io/>
2. **Anthropic MCP Announcement:** <https://www.anthropic.com/news/model-context-protocol>
3. **Web Crawling Research (focused crawling, deduplication):**
 - Chakrabarti, C., Van den Berg, M., & Dom, B. (1999). "Focused crawling using context graphs". VLDB.

- Cho, J., Garcia-Molina, H., & Page, L. (1998). "Efficient crawling through URL ordering". WWW.

4. Distributed Systems:

- Lamport, L. (1978). "Time, Clocks, and the Ordering of Events in a Distributed System". Communications of the ACM.
- Maymounkov, P., & Mazières, D. (2002). "Kademlia: A peer-to-peer information system based on XOR metric". IPSN.

5. Web Compliance & Standards:

- Robots Exclusion Protocol: <https://www.robotstxt.org/>
- Google Search Essentials: <https://developers.google.com/search/docs>

End of Document