

Contents

Razionale delle Decisioni: Perché Questa Architettura	1
Panoramica	1
Meta-Livello: Perché Questa Architettura Esiste	1
Domanda: Perché non un pattern più semplice?	1
Layer 1: Cognitive Layer	2
Decisione 1.1: Separare Goal Analysis da Planning	2
Decisione 1.2: Planning Strategy - HTN vs Flat Decomposition	3
Decisione 1.3: Reflection Module - Post-Execution vs Real-Time	3
Layer 2: Memory System	4
Decisione 2.1: Tre Tipi di Memory vs Single Unified Memory	4
Decisione 2.2: Episodic Memory - Vector DB vs Relational DB	5
Decisione 2.3: Pattern Cache - Solo Pattern Validati vs Tutti i Pattern	5
Layer 3: Capability Layer	6
Decisione 3.1: Tool Registry - Dynamic Discovery vs Static Configuration	6
Decisione 3.2: Model Router - Tre Tier vs Due Tier	6
Decisione 3.3: Safety Verifier - Validazione Pre E Post vs Solo Pre	7
Decisioni Trasversali	8
Decisione 4.1: Bounded Emergence vs Pure Emergence vs Full Control	8
Decisione 4.2: Synchronous Reflection vs Async Reflection	9
Decisione 4.3: Episodic Memory - Memorizzare Tutti gli Episodi vs Sampling	9
Scelta Pattern Architetturale	10
Decisione 5.1: Perché “Reflective Agent” Pattern vs Alternative	10
Razionale Target Performance	10
Decisione 6.1: Perché Target Success Rate 85-95%	10
Decisione 6.2: Perché Target Costo \$0.05-0.30 Per Task	11
Conclusione: Filosofia Architetturale	11

Razionale delle Decisioni: Perché Questa Architettura

Panoramica

Questo documento spiega il **reasoning** dietro ogni scelta architettonica importante. Per ogni componente e design decision, rispondiamo: **Perché così? Quali alternative? Quali trade-off?**

Meta-Livello: Perché Questa Architettura Esiste

Domanda: Perché non un pattern più semplice?

Decisione: Architettura con 4 layer (Cognitive, Memory, Capability, Infrastructure) e 8+ componenti core.

Alternative Considerate: 1. **Minimal Loop** (3 componenti: LLM + Tools + Memory)
2. **ReAct-style** (Reasoning + Acting in loop semplice)

Razionale:

L'architettura minimal funziona per: - [OK] Task semplici, single-step - [OK] Problem class A (Simple/Variable) - [OK] Best-effort acceptable

Ma **fallisce** per: - [NO] Complex planning (non decomponere efficientemente) - [NO] Learning from experience (no reflection systematic) - [NO] Error recovery sofisticato (no pattern riutilizzabili) - [NO] Safety bounds (no enforcement multi-layer)

Target: Problem Class B (Complex Planning) che rappresenta **70-80% dei casi d'uso pratici** in production.

Trade-off Accettati: - [=] Maggiore complessità iniziale - [=] Overhead di 15-30s per planning/reflection - [=] Più componenti da implementare e mantenere

Benefici: - [OK] 85-95% success rate vs 70-85% per minimal loop - [OK] Learning nel tempo (performance migliora) - [OK] Robust error handling - [OK] Bounded emergence (safe + flexible)

Layer 1: Cognitive Layer

Decisione 1.1: Separare Goal Analysis da Planning

Domanda: Perché non combinare goal analysis e planning in un singolo step?

Decisione: Goal Analysis è componente separato che produce GoalStructure consumato da Planning Engine.

Alternativa: Planning Engine riceve direttamente task natural language e fa parsing + planning insieme.

Razionale:

Separation of Concerns: - Goal Analysis: **Cosa** vogliamo (obiettivi, vincoli, contesto)
- Planning: **Come** ottenerlo (decomposizione, sequencing)

Vantaggi separazione: 1. **Riutilizzabilità:** Stesso goal può richiedere planning diversi in context diversi 2. **Caching:** GoalStructure può essere cachata per task simili 3. **Human-in-Loop:** Umano può validare/modificare goal prima di planning 4.

Chiarezza: Reasoning più trasparente (prima "cosa", poi "come")

Esempio:

Task: "Refactor authentication to use JWT"

Goal Analysis Output:

- Primary Goal: Implementare JWT auth
- Constraints: Mantenere retrocompatibilità, nessun downtime
- Context: App Python/FastAPI, auth esistente basata su sessioni

-> Planning può ora ottimizzare per questi constraint specifici
-> Se constraint cambiano, replanning facile senza re-analysis

Trade-off: +15-30s overhead per analysis separata, ma migliore quality di plan.

Decisione 1.2: Planning Strategy - HTN vs Flat Decomposition

Domanda: Perché Hierarchical Task Network planning invece di flat decomposition?

Decisione: Planning Engine usa HTN-inspired approach con decomposizione ricorsiva multi-level.

Alternative: 1. **Flat Decomposition:** Task -> Lista piatta di subtask atomici 2. **Re-active Planning:** Nessun upfront plan, decide step-by-step

Razionale:

Vantaggi HTN: 1. **Scalabilità:** Gestisce task complessi (50+ step) con deep hierarchies 2. **Astrazione:** Ragionamento a livelli diversi di granularità 3. **Riutilizzabilità:** Subtree di plan riutilizzabili 4. **Raffinamento Incrementale:** Piano può essere raffinato progressivamente

Problemi Flat Decomposition: - [NO] Sovraccarico cognitivo (troppi subtask at once)
- [NO] Difficile gestire dipendenze complesse - [NO] Nessuna astrazione (tutto allo stesso livello)

Problemi Reactive Planning: - [NO] Nessuna previsione (non anticipa problemi) - [NO] Inefficiente (decide ogni volta senza overall strategy)

Esempio:

HTN:

```
Livello 0: Implementare autenticazione JWT
Livello 1: Ricerca | Design | Implementazione | Test | Deploy
Livello 2: Ricerca.1: Trovare librerie
          Ricerca.2: Confrontare opzioni
          Ricerca.3: Scegliere la migliore
          ...
          ...
```

vs Flat:

```
Task -> [Trovare librerie, Confrontare opzioni, Scegliere migliore, Installare lib,
          Configurare, Scrivere codice, Scrivere test, Eseguire test, Deploy, ...]
-> 30+ items allo stesso livello, opprimente
```

Trade-off: HTN più complesso da implementare, ma essential per task complessi.

Decisione 1.3: Reflection Module - Post-Execution vs Real-Time

Domanda: Perché reflection avviene POST-execution invece di real-time durante execution?

Decisione: Reflection è una fase separata dopo task completion.

Alternativa: Continuous reflection durante execution (dopo ogni subtask).

Razionale:

Vantaggi Post-Execution: 1. **Quadro Completo:** Può analizzare intero episodio con outcome finale 2. **Nessun Overhead di Esecuzione:** Non rallenta task execu-

tion 3. **Migliore Estrazione Pattern**: Vede sequenze complete, non frammenti 4. **Asincrono**: Può girare in background senza bloccare utente

Problemi Real-Time: - [NO] Overhead significativo (+30-50% execution time) - [NO] Pattern incompleti (non vede fine task) - [NO] Distrazione (agent deve cambiare context costantemente)

Quando Real-Time È Meglio: Sarebbe preferibile per long-running tasks (>10 min) dove mid-course correction è critica. Ma i target use cases sono task di 2-10 min.

Compromesso: Monitoring real-time (rilevamento stuck/failures) + Reflection post-execution (apprendimento).

Layer 2: Memory System

Decisione 2.1: Tre Tipi di Memory vs Single Unified Memory

Domanda: Perché tre memory systems separati (Working, Episodic, Pattern) invece di un'unica memoria?

Decisione: Tre sottosistemi specializzati con caratteristiche diverse.

Alternativa: Unified memory con single storage + query interface.

Razionale:

Requisiti Differenti:

Tipo Memoria	Dim.	Lifetime	Pattern Accesso
Working	~20KB	Task	Random, freq.
Episodic	GBs	Forever	Semantic search
Pattern	MBs	Forever	Match + rank

Unified memory non può ottimizzare per questi pattern diversi: - Working necessita **velocità in-memory** - Episodic necessita **semantic search** su larga scala - Pattern necessita **matching strutturato** con validazione

Tentativo di Unificazione richiederebbe compromessi: - Accesso lento (disk-based) -> distrugge performance Working Memory - Nessuna specializzazione -> subottimale per ogni use case

Simile alla Gerarchia Cache CPU: Cache L1/L2/L3 vs memoria principale. Dimensioni, velocità, scopi diversi.

Trade-off: Più complessità (3 systems vs 1), ma **10-100x** migliore performance per ogni use case.

Decisione 2.2: Episodic Memory - Vector DB vs Relational DB

Domanda: Perché vector database per Episodic Memory invece di relational DB?

Decisione: Irido: Vector DB (semantic search) + Document DB (structured storage).

Alternativa: Solo relational DB con full-text search.

Razionale:

Vector DB Essenziale per similarità semantica:

Query: "Autenticazione fallita con token non valido"

Relational DB con keyword search:

- > Trova episodi contenenti parole "autenticazione", "fallita", "token"
- > Potrebbe mancare: "errore verifica JWT", "accesso non autorizzato", ecc.
- > Matching per keyword troppo rigido

Vector DB con embeddings:

- > Trova episodi SEMANTICAMENTE simili
- > Trova "errore verifica JWT" (parole diverse, stesso significato)
- > Trova "non autorizzato per credenziali errate" (concetto correlato)

Perché Irido: - Vector DB per **scoperta** (trovare simili) - Document DB per **recupero** (ottenere episodio completo per ID) - Document DB per **query strutturate** (filtrare per attributi)

Trade-off: Più infrastruttura (2 database), ma **essenziale** per qualità del recupero semantico.

Decisione 2.3: Pattern Cache - Solo Pattern Validati vs Tutti i Pattern

Domanda: Perché solo pattern validati in cache, invece di tutti i pattern scoperti?

Decisione: Pattern Cache contiene solo pattern VALIDATI (confidence > threshold, significatività statistica).

Alternativa: Memorizzare tutti i pattern scoperti, lasciare che l'agent filtri.

Razionale:

Qualità su Quantità: - 1000 pattern di bassa qualità -> rumore, confusione - 100 pattern di alta qualità -> utilizzabili, affidabili

Il Processo di Validazione Garantisce: 1. **Significatività Statistica:** Non solo coincidenza 2. **Evidenza Sufficiente:** ≥ 5 episodi di supporto 3. **Consistenza:** Tasso di successo > threshold 4. **Nessuna Contraddizione:** Nessun conflitto con altri pattern validati

Pattern CANDIDATI: Memorizzati separatamente, tracciati fino a evidenza sufficiente, poi promossi o rifiutati.

Perché Non Memorizzare Tutti: L'agent deve decidere rapidamente durante planning. Dover filtrare/validare pattern ogni volta distrugge la latenza.

Trade-off: Alcuni pattern potenzialmente utili ritardati fino a validazione, ma **maggiore affidabilità** dei pattern utilizzati.

Layer 3: Capability Layer

Decisione 3.1: Tool Registry - Dynamic Discovery vs Static Configuration

Domanda: Perché dynamic tool discovery invece di static tool configuration?

Decisione: Tools registrati dinamicamente con capability tags, scoperti a runtime.

Alternativa: Lista hardcoded di tools disponibili, agent conosce in anticipo.

Razionale:

Dynamic Discovery Abilità: 1. **Estensibilità:** Aggiungere nuovi tool senza modifiche al codice agent 2. **Context-Specific:** Ambienti diversi hanno tool diversi 3.

Basato su Permessi: Mostrare solo tool per cui l'utente ha permesso 4. **Versioning:** Multiple versioni dello stesso tool coesistono

Esempio:

Configurazione Statica:

```
Codice Agent: tools = [read_file, write_file, http_get, ...]  
-> Per aggiungere nuovo tool: modificare codice agent, rideploy
```

Dynamic Discovery:

```
Agent: "Ho bisogno di leggere file JSON"  
Registry: "Ecco i tool con capability 'read_file': [A, B, C]"  
Agent: "B sembra il migliore, lo uso"  
-> Per aggiungere nuovo tool: registrare in registry, disponibilità immediata
```

Trade-off: ~10-20ms overhead per discovery, ma la flessibilità vale per sistema estensibile.

Decisione 3.2: Model Router - Tre Tier vs Due Tier

Domanda: Perché three model tiers (Small/Medium/Large) invece di two (Small/Large)?

Decisione: Sistema a tre tier con profili costo/capability distinti.

Alternativa: Due tier (fast/cheap vs slow/expensive).

Razionale:

Medium Tier è “Sweet Spot” per la maggior parte dei task:

Sistema a Due Tier:

```
Small: $0.001/1K, buono per <20% task  
Large: $0.05/1K, necessario per >80% task  
Costo medio: ~$0.04/1K
```

Sistema a Tre Tier:

Small: \$0.001/1K, buono per ~20% task

Medium: \$0.015/1K, buono per ~60% task

Large: \$0.05/1K, necessario per ~20% task

Costo medio: ~\$0.018/1K

Riduzione costi di 2.2x indirizzando la maggioranza a Medium tier.

Perché Medium Tier Funziona: - I modelli medium moderni (GPT-4o-mini, Claude Sonnet) sono **sufficientemente capaci** per la maggior parte di coding/analysis - Solo problemi veramente nuovi/complessi richiedono i modelli più grandi - I modelli small sono veramente troppo limitati per la maggior parte del lavoro reale

Trade-off: Maggiore complessità di routing, ma **risparmi significativi** (~50-70%).

Decisione 3.3: Safety Verifier - Validazione Pre E Post vs Solo Pre

Domanda: Perché validation sia BEFORE che AFTER action execution?

Decisione: Validazione multi-layer: Input -> Action Authorization -> Execution -> Output Validation.

Alternativa: Solo pre-execution validation (autorizzare azione, eseguire, restituire risultato).

Razionale:

Pre-Validation Non Sufficiente:

Esempio:

Action: "Leggere file user_data.json"

Pre-validation: [v] Ha permesso lettura, path valido

Execution: Legge file, restituisce:

```
{  
  "users": [...],  
  "api_key": "sk-abc123xyz...", <- Leak!  
  "admin_password": "secret123" <- Leak!  
}
```

Senza Post-Validation: Restituisce questo all'agent -> VIOLAZIONE SICUREZZA

Con Post-Validation: Rileva segreti, RIFIUTA output

Defense in Depth: 1. Pre-validation: Prevenire azioni dannose 2. Post-validation: Catturare output dannosi (anche da azioni "sicure")

L'Output Può Essere Pericoloso Anche Se L'Azione È Sicura: - Segreti trapelati (chiavi API, password) - Codice malevolo generato - PII esposti - Violazioni di policy

Trade-off: +50-100ms per azione per validazione output, ma **essenziale** per sicurezza.

Decisioni Trasversali

Decisione 4.1: Bounded Emergence vs Pure Emergence vs Full Control

Domanda: Perché bounded emergence invece di full emergence o full control?

Decisione: Approccio ibrido - reasoning LLM entro bounds esplicativi.

Alternative: 1. **Pure Emergence:** LLM decide tutto, nessun guardrail 2. **Full Control:** Tutta logica esplicita, nessun reasoning LLM

Razionale:

Problemi Pure Emergence: - [NO] Imprevedibile (può fare qualsiasi cosa) - [NO] Non sicuro (nessuna garanzia) - [NO] Non-deterministico (difficile da debuggare)

Problemi Full Control: - [NO] Rígido (fallisce su input inattesi) - [NO] Overhead ingeneristico (codificare ogni behavior) - [NO] Non generalizza

Bounded Emergence = Il Meglio di Entrambi:

BOUNDS (Esplicativi):

- Validazione input (schema, rilevamento injection)
- Whitelist tool (operazioni consentite)
- Verifica output (formato, sicurezza)
- Limiti risorse (tempo, costo, token)

EMERGENCE (LLM):

- Decomposizione goal
- Selezione strategia
- Approcci error recovery
- Riconoscimento pattern

Risultato: Flessibilità DENTRO bounds di sicurezza.

Esempio:

Task: "Ottimizzare query database"

Bounds Prevengono:

- Accesso diretto database (non whitelisted)
- Eliminazione tabelle
- Loop infiniti

Emergence Abilita:

- Analizzare piano query
- Suggerire miglioramenti indici
- Riscrivere struttura query
- Provare approcci multipli

-> Esplorazione sicura dello spazio soluzioni

Trade-off: Bounds richiedono design attento, ma **essenziali** per sicurezza produzione.

Decisione 4.2: Synchronous Reflection vs Async Reflection

Domanda: Perché reflection asincrona invece di bloccante?

Decisione: Reflection gira asincronamente dopo task completion, non blocca response a user.

Alternativa: Synchronous reflection - completare reflection prima di restituire risultato.

Razionale:

User Experience:

Synchronous:

Task completa → Reflection (20–40s) → Restituisce risultato

Tempo attesa utente: Tempo task + 20–40s

Asynchronous:

Task completa → Restituisce risultato immediatamente

Reflection in background

Tempo attesa utente: Solo tempo task

Reflection Non Time-Critical: - Apprendimento per **task futuri**, non quello corrente
- L'utente non deve attendere estrazione pattern - Può essere batch/ottimizzato di notte

Quando Synchronous È Meglio: Se i risultati reflection servono immediatamente per validazione. Ma questo è **Verification** (synchronous), non Reflection (apprendimento).

Trade-off: Leggero ritardo prima che l'apprendimento sia disponibile, ma **UX molto migliore**.

Decisione 4.3: Episodic Memory - Memorizzare Tutti gli Episodi vs Sampling

Domanda: Perché memorizzare TUTTI gli episodi invece di solo un subset rappresentativo?

Decisione: Memorizzare ogni episodio (con eventuale consolidamento/archiviazione).

Alternativa: Memorizzare solo episodi di successo, o campione casuale.

Razionale:

I Fallimenti Sono Preziosi: - Apprendimento da pattern di fallimento - Identificare cosa NON fare - Debug problemi sistematici

La Long Tail Conta: - Situazioni rare potrebbero essere importanti - Sampling potrebbe perdere casi limite critici - Storage è economico, insights sono costosi

Consolidamento Gestisce Scala: - Episodi simili uniti nel tempo - Vecchi episodi archiviati - Ma inizialmente, conservare tutto

Quando Sampling È Meglio: Se ci sono preoccupazioni privacy (non si vogliono memorizzare tutti i dati utente). Ma si può anonomizzare invece.

Trade-off: Costo storage (~\$10-200/mese a seconda della scala), ma **apprendimento completo**.

Scelta Pattern Architetturale

Decisione 5.1: Perché “Reflective Agent” Pattern vs Alternative

Domanda: Perché scegliere Reflective Agent pattern (Pattern 2) come reference architecture?

Alternative dal Framework: 1. **Minimal Loop** (Pattern 1) 2. **Verified Agent** (Pattern 3) 3. **Reactive Agent** (Pattern 4) 4. **Multi-Agent** (Pattern 5)

Matrice Decisionale:

Pattern	Successo	Latenza	Costo	Casi Uso
Minimal	70-85%	<5s	Basso	20%
Reflective	85-95%	5-30s	Medio	70%
Verified	95-99%	Variab.	Alto	5%
Reactive	80-90%	<2s	Medio	3%
Multi-Agent	90-95%	10-60s	Alto	2%

Reflective Agent Vince perché: 1. **Copertura 70-80%**: Maggioranza dei casi d'uso pratici 2. **Trade-off Bilanciati**: Non troppo semplice, non troppo complesso 3. **Capacità di Apprendimento**: Migliora nel tempo 4. **Production-ready**: Adatto per deployment reale

Quando Altri Sono Meglio: - Minimal: Q&A semplice, prototipazione - Verified: Safety-critical (medico, finanziario) - Reactive: Controllo real-time (<100ms) - Multi-Agent: Scenari collaborazione complessa

Reflective è “Goldilocks”: Non troppo semplice, non troppo complesso, giusto per la maggior parte.

Razionale Target Performance

Decisione 6.1: Perché Target Success Rate 85-95%

Domanda: Perché target 85-95% invece di 99%+?

Razionale:

Rendimenti Decrescenti:

70% -> 85%: Raggiungibile con buona architettura
85% -> 95%: Richiede reflection + apprendimento
95% -> 99%: Richiede verifica estensiva (costo 10x)
99% -> 99.9%: Richiede metodi formali (costo 100x)

Analisi Casi d'Uso: - 70% non production-ready (troppi fallimenti) - 85-95% accettabile per la maggior parte automazione - 99%+ necessario solo per sistemi critici

Baseline Performance Umana: Gli sviluppatori hanno ~85-90% tasso successo al primo tentativo. Eguagliare l'umano è obiettivo ragionevole.

Decisione 6.2: Perché Target Costo \$0.05-0.30 Per Task

Razionale:

Sostenibilità Economica:

Complessità task: 5-15 minuti di tempo sviluppatore umano

Costo umano: \$50/ora -> \$4-12 per task

Target costo agent: \$0.05-0.30 -> 10-100x più economico

-> ROI convincente anche a \$0.30/task

Fattibilità Tecnica: - ~60-125K token per task complesso - Mix di modelli small/medium/large - Costo medio per token: \$0.005-0.015/1K - Matematica: $80K \text{ token} * \$0.010/1K = \0.80 (caso peggiore) - Con routing: ~\$0.15-0.25 tipico

Margine per Futuro: Target costi assumono che i modelli diventino più economici nel tempo (storicamente vero).

Conclusione: Filosofia Architetturale

Questa architettura riflette una **filosofia di design pragmatico**:

1. **Non la Più Semplice:** Deliberatamente più complessa di minimal loop
2. **Non la Più Complessa:** Deliberatamente meno complessa di architettura cognitiva completa
3. **Bilanciata:** Sweet spot per uso produzione

Principi Guida: - [OK] **Empirica:** Decisioni basate su dati, non ideologia - [OK] **Pragmatica:** Scegliere cosa funziona, non cosa è elegante - [OK] **Evolvibile:** Può iniziare semplice, aggiungere complessità quando necessario - [OK] **Onesta:** Riconoscere trade-off, nessuna soluzione magica

Non Universale: Questa è ottimale per ~70-80% casi d'uso. Altri casi d'uso richiedono architetture diverse. Il Framework (repo principale) fornisce guidance.

Validazione Necessaria: Queste sono decisioni informate da analisi, ma **richiedono validazione empirica** attraverso implementazione e testing.

Fine Razionale Decisioni

Per guidance implementativa, vedere specifiche componenti individuali. Per architecture alternative, vedere `/03-architecture-patterns.md` nel framework principale.