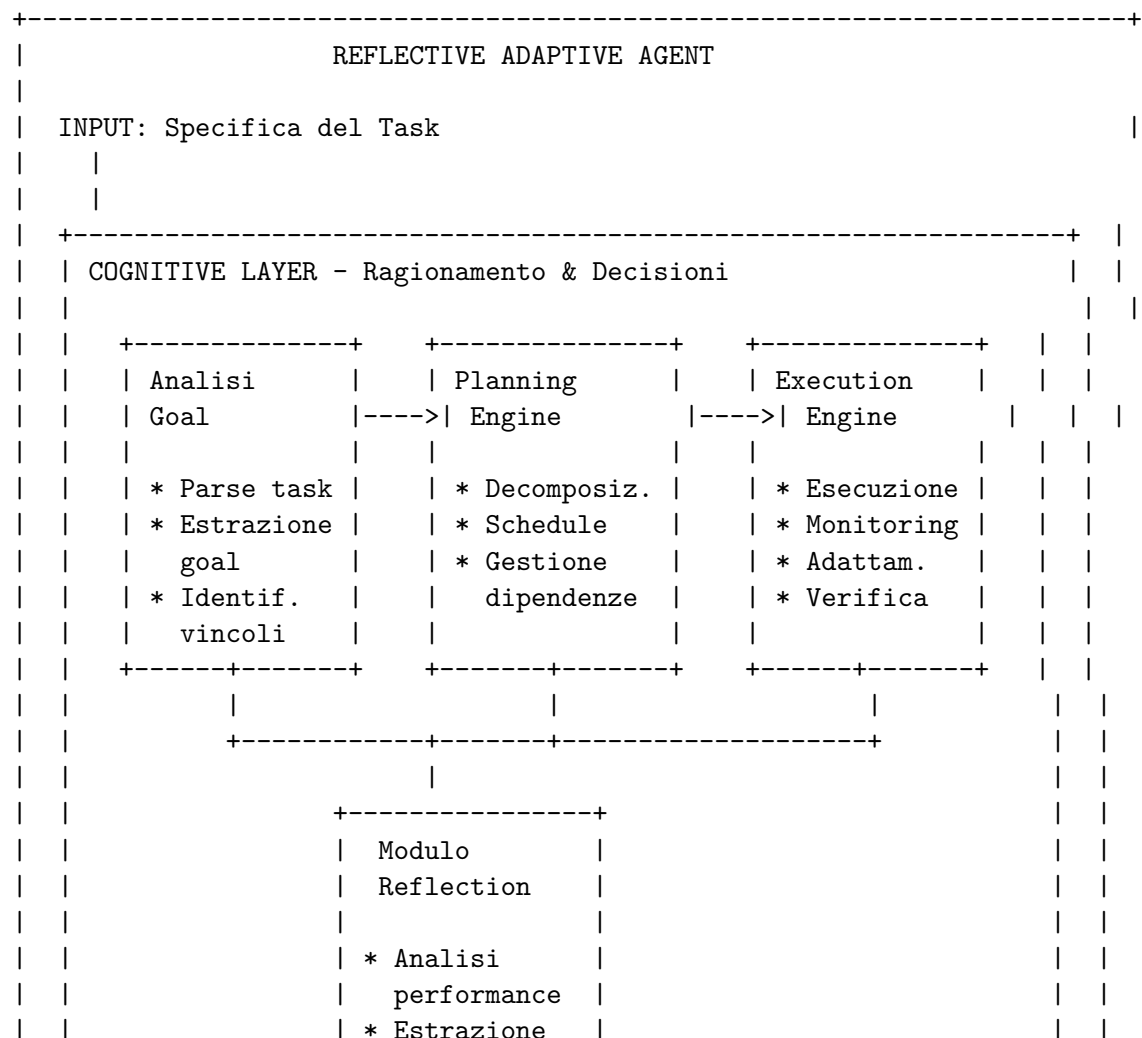


Contents

Architettura di Sistema - Specifiche Complete		1
Diagramma Generale		1
Flusso di Esecuzione - Loop Principale		3
Interazione tra Componenti - Diagramma di Sequenza		4
Architettura Flussi Dati		5
Macchina a Stati - Stati di Esecuzione		7
Dipendenze tra Componenti		8
Architettura di Deployment		9
Caratteristiche di Performance		10
Breakdown Budget Latenza		10
Caratteristiche di Scalabilità		11

Architettura di Sistema - Specifiche Complete

Diagramma Generale



		pattern			
		* Aggiornamento			
		strategie			
		+-----+			
		+-----+			
		+-----+			
		MEMORY LAYER - Gestione Stato & Conoscenza			
		+-----+	+-----+	+-----+	
		Working	Episodic	Pattern	
		Memory	Memory	Cache	
		* Contesto	* Episodi	* Strategie	
		corrente	passati	apprese	
		* Variabili	* Risultati	* Pattern	
		attive	* Embeddings	successo	
		* Stato temp	* Retrieval	* Euristiche	
			similitudine		
		+-----+	+-----+	+-----+	
		+-----+	+-----+	+-----+	
		+-----+	+-----+	+-----+	
		+-----+	+-----+	+-----+	
		CAPABILITY LAYER - Azioni Esterne & Intelligenza			
		+-----+	+-----+	+-----+	
		Tool	Model	Safety	
		Registry	Router	Verifier	
		* Discovery	* Routing al	* Validazione	
		* Binding	modello	input	
		* Esecuzione	appropriato	* Controllo	
		* Validazione	* Ottimizz.	bounds	
			costi	* Enforce	
				safety	
		+-----+	+-----+	+-----+	
		+-----+	+-----+	+-----+	
		+-----+	+-----+	+-----+	
		+-----+	+-----+	+-----+	
		INFRASTRUCTURE LAYER - Observability & Gestione Risorse			

		+-----+	+-----+	+-----+		
		Observability	Resource	Error		
		System	Manager	Handler		
		* Tracing	* Tracking	* Rilevazione		
		* Metriche	budget	* Classific.		
		* Logging	* Rate	* Recovery		
		* Monitoring	limiting	* Escalation		
			* Throttling			
		+-----+	+-----+	+-----+		
	+-----+					
	OUTPUT: Risultato Task + Conoscenza Aggiornata					
+-----+						

Flusso di Esecuzione - Loop Principale

INIZIO

```

|
+--> [1. RICEZIONE TASK]
|
|   +- Parse input
|   +- Validazione formato
|   +- Estrazione parametri
|
+--> [2. ANALISI GOAL]
|
|   +- Identificazione obiettivi
|   +- Estrazione vincoli
|   +- Definizione criteri successo
|   +- Classificazione complessità
|
+--> [3. RECUPERO CONTESTO]
|
|   +- Query episodic memory (task simili passati)
|   +- Controllo pattern cache (strategie apprese)
|   +- Caricamento working memory (sessione corrente)
|
+--> [4. PIANIFICAZIONE]
|
|   +- Generazione piano alto livello
|   +- Decomposizione in subtask
|   +- Identificazione dipendenze
|   +- Stima risorse (token, tempo, costo)
|   +- Selezione strategia esecuzione
|
+--> [5. LOOP ESECUZIONE]
```

```

|
| +- PER ogni step nel piano:
| | |
| | | +--> [5a. Pre-esecuzione]
| | | | +- Verifica sicurezza
| | | | +- Controllo risorse
| | | | +- Decisione routing modello
| | | |
| | | +--> [5b. Esecuzione]
| | | | +- Reasoning LLM (se necessario)
| | | | +- Esecuzione tool (se necessario)
| | | | +- Cattura output
| | | |
| | | +--> [5c. Post-esecuzione]
| | | | +- Verifica risultato
| | | | +- Controllo criteri successo
| | | | +- Aggiornamento working memory
| | | |
| | | +--> [5d. Adattamento]
| | | | +- Se successo: Continua
| | | | +- Se fallimento: Strategia recovery
| | | | +- Se bloccato: Ripianificazione
| | |
| +--> [6. VERIFICA]
| | |
| | | +- Validazione output finale
| | | +- Controllo criteri successo
| | | +- Verifica sicurezza finale
| | |
| +--> [7. REFLECTION]
| | |
| | | +- Analisi performance episodio
| | | +- Estrazione pattern di successo
| | | +- Identificazione miglioramenti
| | | +- Aggiornamento pattern cache
| | | +- Memorizzazione in episodic memory
| | |
| +--> [8. RESTITUZIONE RISULTATO]

```

Interazione tra Componenti - Diagramma di Sequenza

User	Analisi Goal	Planning Engine	Memory System	Model Router	Tool Registry	Safety Verifier
+--task-->						
	+goals-->					

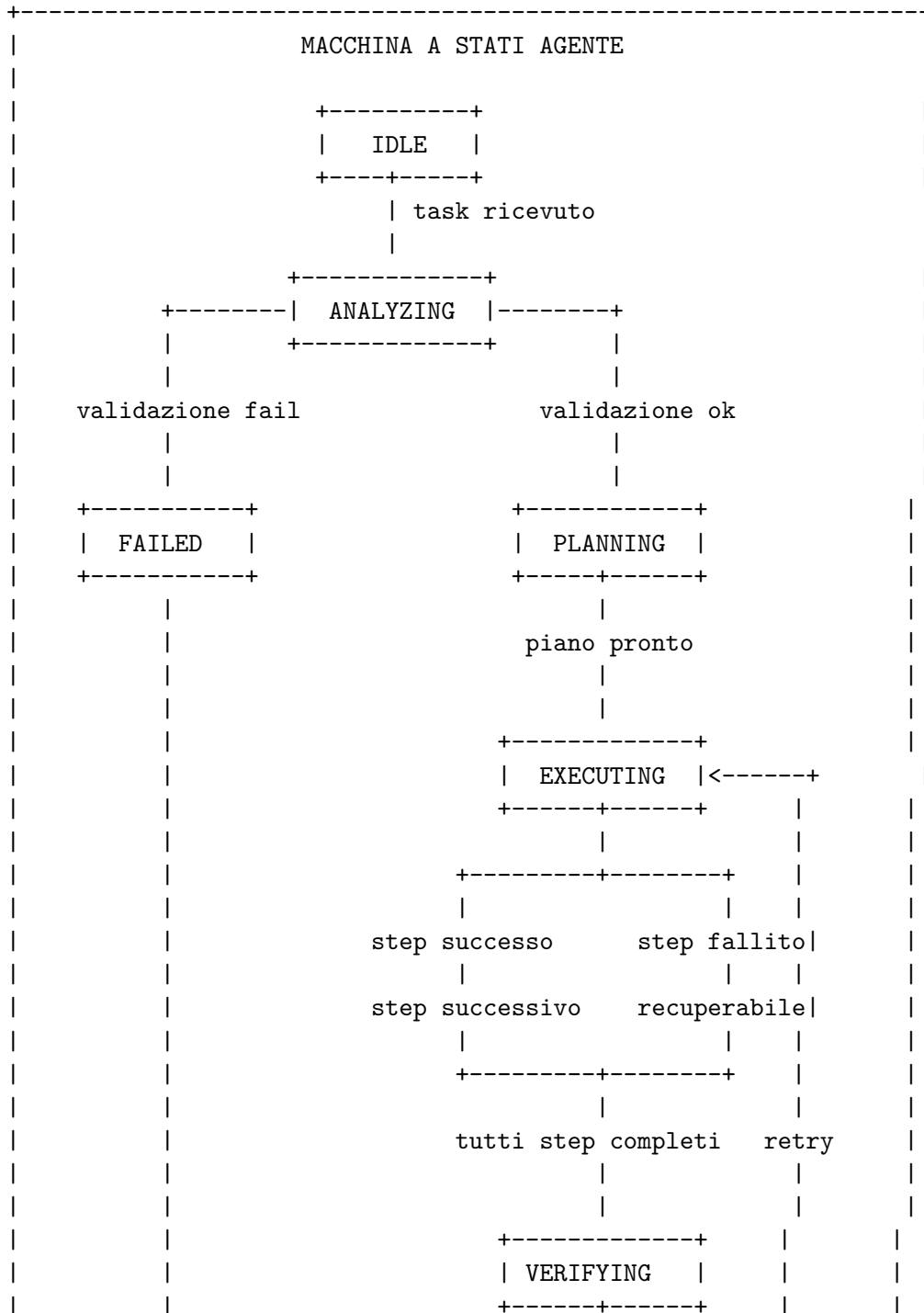


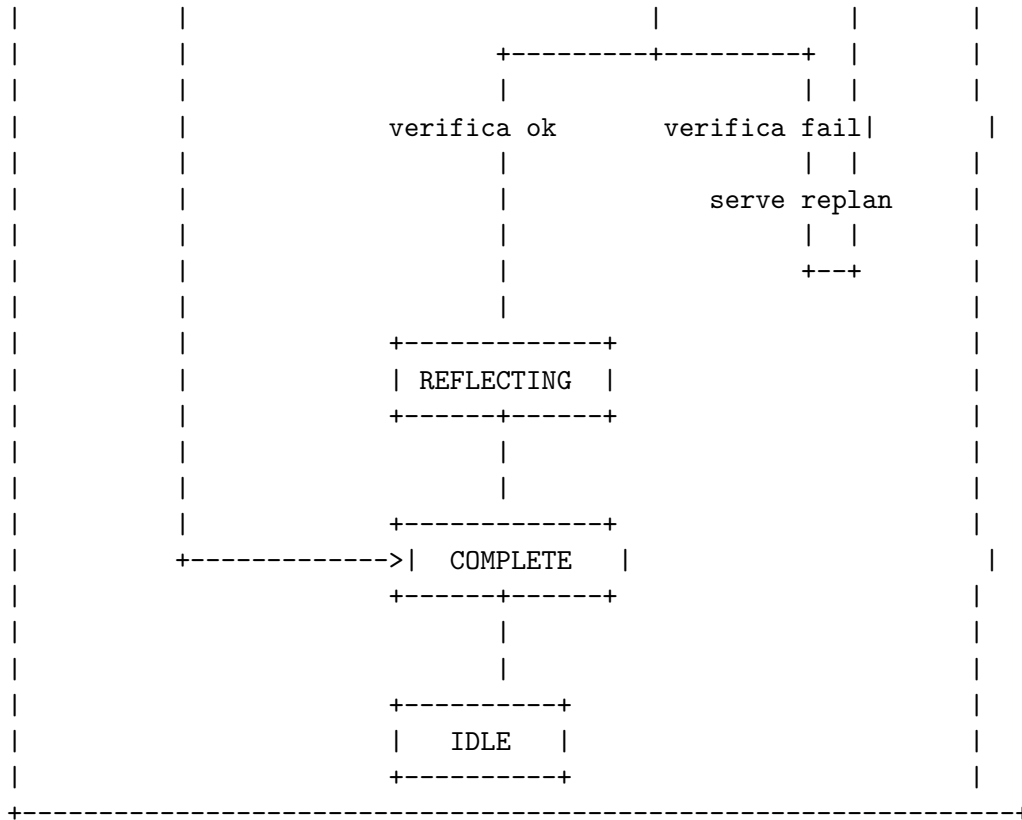
```

| | OUTPUT | |
| | RISULTATO | |
| +-----+ |
+-----+

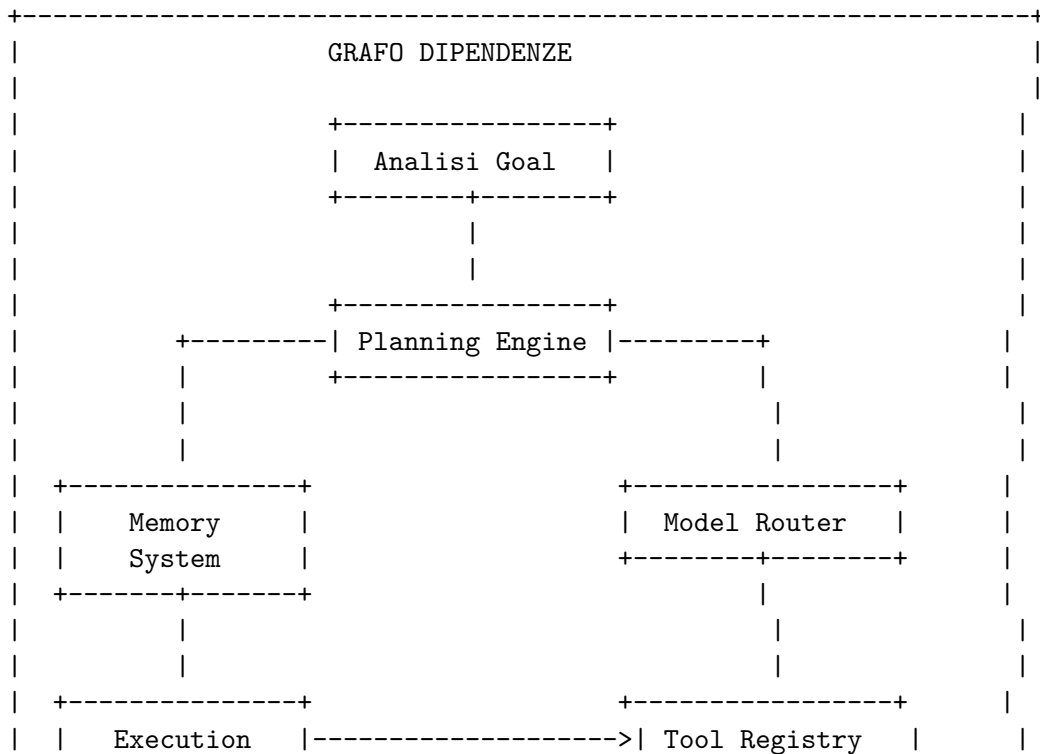
```

Macchina a Stati - Stati di Esecuzione





Dipendenze tra Componenti




```

+-----+
|
+-----+
| TIER CACHE |
|
| +-----+
| | Redis / Memcached |
| | (Pattern Cache, Stato Sessione) |
| +-----+
+-----+
|
+-----+
| TIER DATI |
|
| +-----+ +-----+ +-----+ |
| | PostgreSQL | | Pinecone | | S3/Blob | |
| | (Metadata) | | (Episodic | | (Log, | |
| | | | Memory) | | Artifact) | |
| +-----+ +-----+ +-----+ |
+-----+
|
+-----+
| SERVIZI ESTERNI |
|
| +-----+ +-----+ +-----+ |
| | Anthropic | | OpenAI | | Tool | |
| | Claude | | GPT-4 | | Custom | |
| +-----+ +-----+ +-----+ |
+-----+
|
+-----+
| OBSERVABILITY |
|
| +-----+ +-----+ +-----+ |
| | Prometheus | | Grafana | | ELK | |
| | (Metriche) | | (Dashboard) | | (Log) | |
| +-----+ +-----+ +-----+ |
+-----+

```

+-----+			
Analisi Goal	0.5s	2%	
Recupero Contesto	1.0s	5%	
Pianificazione	2.0s	10%	
Esecuzione:			
- Chiamate model	15.0s	70%	
- Esecuzione tool	3.0s	13%	
Verifica	0.5s	2%	
Reflection	1.0s	5%	
+-----+			
TOTALE	~23s	~100%	
+-----+			

Strategie di ottimizzazione:

- Routing modelli (usa modelli veloci quando sufficiente)
- Esecuzione parallela tool
- Riutilizzo pattern cached
- Streaming per latenza percepita

Caratteristiche di Scalabilità

+-----+			
Dimensione	Corrente	Target	Strategia Scaling
+-----+			
Utenti Concurrent	10-50	100-500	Orizzontale
Task/Ora	100-500	5K-10K	Orizzontale
Episodi Memory	10K	1M+	Scala Vector DB
Pattern Cache	1K	10K+	Cache distribuita
Latenza Media	5-30s	<30s	Ottimizz. modelli
Success Rate	85-95%	>90%	Apprendimento cont
+-----+			

Prossimo: 02-cognitive-layer.md -> Specifiche dettagliate dei componenti di ragionamento