

GUIDA OPERATIVA · FEBBRAIO 2026

Il Costo Reale dell'AI

PER IMPRENDITORI E MANAGER

Quanto costa davvero l'AI nella tua azienda? E quando vale la pena investirci? Qui trovi i numeri reali — per funzione, per task, per livello di modello — e un metodo semplice per decidere: Go, Watch o Stop. Niente teoria, solo dati operativi.

Budget AI

ROI operativo

Governance AI

PMI

Come usare questo documento: identifica i task AI del tuo team → assegna ogni task a un livello (A/B/C/D) → stima volume e costo mensile con la formula in Sezione 5 → definisci il KPI minimo → ogni mese applica la soglia Go/Watch/Stop.

Sezione 1 – I 4 livelli di AI

Dimentica i nomi tecnici dei modelli. Esistono 4 livelli operativi:

LIVELLO	NOME BUSINESS	QUANDO USARLO	COSTO/TASK (AGENTICO)	COSTO/TASK (UNA TANTUM)	MODELLI DI RIFERIMENTO
A	Consulente Senior AI	Analisi ad alto impatto, decisioni critiche, documenti strategici	€0,05 – €0,30	€0,02 – €0,08	Claude Opus 4.5/4.6 · GPT-4o
B	Specialista AI	Lavoro standard di qualità: report, bozze qualificate, analisi correnti	€0,01 – €0,08	€0,005 – €0,02	Claude Sonnet 4.5 · Gemini 2.5 Pro · Mistral Large
C	Operativo AI	Volume, ripetizioni, classificazione, riassunti standard	€0,001 – €0,015	€0,0005 – €0,005	Claude Haiku 4.5 · GPT-4o-mini · Gemini 2.5 Flash
D	Automazione	Task deterministici e prevedibili, workflow fissi	€0,00 – €0,003	<€0,001	Mistral Nemo · modelli locali (Ollama)

¹ **Uso agentic:** l'AI esegue più passaggi in autonomia (3-8 chiamate API per task).

² **Una tantum:** singola chiamata API, supervisione umana diretta.

-67%

CALO PREZZI MODELLI
PREMIUM VS 2024

20/80

LIVELLO A: MAX 20% DEI TASK,
80% DEL VALORE

-90%

RISPARMIO MASSIMO CON
BATCH + CACHING

Regola del 20/80: il Livello A dovrebbe coprire al massimo il 20% dei task, ma genera l'80% del valore. Usarlo su tutto è la causa principale di costi fuori controllo.

Sezione 2 – Quanto costa un agente AI?

SCENARIO	COSTO GIORNALIERO	COSTO ANNUO EQUIVALENTE	CONFRONTO RISORSE UMANE	PUNTO DI PAREGGIO
Agente leggero – Livello C/D	€3 – €20	€1.100 – €7.300	< costo di un'ora di lavoro/giorno	Quasi sempre conveniente
Agente standard – Livello B	€20 – €80	€7.300 – €29.200	~costo di 1-3 ore lavorative/giorno	Conveniente se automatizza ≥ 4 ore/giorno
Agente premium NON ottimizzato – Livello A su tutto	€100 – €300	€36.000 – €108.000	Comparabile a un junior o part-time	A rischio: serve ROI documentato ≥ 3x
Agente premium OTTIMIZZATO – Livello A solo dove serve	€15 – €50	€5.500 – €18.200	~costo di 30-90 min lavoro/giorno	Go se automatizza ≥ 2 ore/giorno qualificato

Il range €100-300/giorno riguarda agenti mal configurati o con contesti molto lunghi non ottimizzati. Il problema non è il modello: è usare il Livello A su task da Livello C, senza batch né caching.

Sezione 3 – Matrice per funzione aziendale

3.1 MARKETING & COMUNICAZIONE

TASK	FREQUENZA	LIVELLO AI	COSTO/TASK	VOLUME MENSILE	COSTO MENSILE STIMATO	KPI MINIMO	SOGLIA
Bozze post social / newsletter	Alta	B	€0,03 – 0,08	60 bozze	€2 – €5	≥ 3 bozze utilizzabili/ora; ≥ 50% tempo risparmiato	Stop se < 30% approvazione
Variazioni copy campagna (A/B test)	Media	C	€0,005 – 0,02	100 variazioni	€0,50 – €2	≥ 5 varianti corrette/set; zero errori di brand	Stop se copy fuori brand
Analisi competitor + positioning	Bassa	A→B	€0,20 – 0,40	4 report	€0,80 – €1,60	1 report completo, tutte le sezioni, ≤ 30 min produzione	Watch se > 45 min
Risposta commenti / community	Alta	C	€0,002 – 0,01	300 risposte	€0,60 – €3	≥ 90% risposte in tono, zero escalation per errore AI	Stop se > 2% errori tono
Produzione brief creativo	Bassa	B	€0,05 – 0,15	8 brief	€0,40 – €1,20	Brief completo, zero rework strutturale	Watch se > 1 ciclo revisione

3.2 COMMERCIALE & VENDITE

TASK	FREQUENZA	LIVELLO AI	COSTO/TASK	VOLUME MENSILE	COSTO MENSILE STIMATO	KPI MINIMO	SOGLIA
Offerte commerciali personalizzate	Media	B	€0,05 – 0,15	30 offerte	€1,50 – €4,50	Bozza pronta in < 10 min; ≥ 80% usata senza riscrittura	Watch se riscrittura > 30%
Outreach email a freddo	Alta	B→C	€0,01 – 0,05	200 email	€2 – €10	≥ 25% open rate; ≥ 5% risposta	Stop se < 15% open rate
Riassunto CRM / note post-call	Alta	C	€0,003 – 0,01	80 note	€0,24 – €0,80	≥ 95% action item corretti; zero task critici persi	Stop se > 1 task critico perso/mese
Analisi pipeline e forecast	Bassa	A	€0,15 – 0,40	2 sessioni	€0,30 – €0,80	≥ 1 insight azionabile; previsione entro ±15%	Watch se insight non azionabili
Preparazione negoziazione (brief)	Bassa	A	€0,20 – 0,50	5 brief	€1 – €2,50	Brief completo in < 20 min; ≥ 1 leva identificata	Watch se leva già note

3.3 OPERAZIONI & AMMINISTRAZIONE

TASK	FREQUENZA	LIVELLO AI	COSTO/TASK	VOLUME MENSILE	COSTO MENSILE STIMATO	KPI MINIMO	SOGLIA
Sintesi riunione + action item	Alta	C	€0,005 – 0,015	40 riunioni	€0,20 – €0,60	≥ 95% action item corretti; distribuzione entro 30 min	Stop se > 1 action item critico mancante/mese
Report operativi ricorrenti	Alta	C→B	€0,01 – 0,05	20 report	€0,20 – €1	≤ 5% errori formato; ≥ 60% riduzione tempo produzione	Stop se formato instabile
Revisione contratti standard	Bassa	A	€0,05 – 0,25	10 documenti	€0,50 – €2,50	Tutte le clausole critiche verificate	Stop: non sostituisce legale
Classificazione e routing documenti	Alta	D	€0,001 – 0,005	500 documenti	€0,50 – €2,50	≥ 98% routing corretto	Stop se < 95%
Risposta email standard / FAQ interne	Alta	C	€0,002 – 0,008	150 email	€0,30 – €1,20	≥ 85% risposte inviate senza intervento umano	Watch se escalation > 15%

3.4 HR & FORMAZIONE

TASK	FREQUENZA	LIVELLO AI	COSTO/TASK	VOLUME MENSILE	COSTO MENSILE STIMATO	KPI MINIMO	SOGLIA
Redazione job description	Media	B	€0,04 – 0,12	8 JD	€0,32 – €0,96	JD pronta in < 15 min; ≥ 1 revisione max	Stop se sempre 2+ revisioni
Pre-screening CV (classificazione)	Alta	C	€0,003 – 0,01	200 CV	€0,60 – €2	≥ 90% classificazione corretta; zero falsi negativi critici	Stop se > 5% falsi negativi
Materiali onboarding / FAQ dipendenti	Bassa	B	€0,05 – 0,15	10 documenti	€0,50 – €1,50	Documenti usabili senza editing; ≥ 70% riduzione tempo	Watch se editing > 30%
Sintesi feedback / survey interni	Bassa	B	€0,08 – 0,20	4 report	€0,32 – €0,80	Pattern principali identificati; ≥ 1 insight azionabile	Watch se nessun insight nuovo
Piano formazione personalizzato	Bassa	A	€0,20 – 0,50	5 piani	€1 – €2,50	Piano completo con obiettivi e timeline; approvato al primo round	Watch se 2+ revisioni

3.5 MANAGEMENT & STRATEGIA

TASK	FREQUENZA	LIVELLO AI	COSTO/TASK	VOLUME MENSILE	COSTO MENSILE STIMATO	KPI MINIMO	SOGLIA
Analisi scenario / decisione pricing	Bassa	A	€0,05 – 0,20	2 sessioni	€0,10 – €0,40	≥ 1 decisione azionabile; chiarezza > baseline	Stop se output già noto
Deck presentazione board / investitori	Bassa	A	€0,05 – 0,25	2 deck	€0,10 – €0,50	Struttura logica; ≤ 1 ciclo di revisione strutturale	Watch se struttura sempre rifatta
Monitoraggio stampa / segnali mercato	Alta	C→B	€0,01 – 0,05	60 digest	€0,60 – €3	≥ 3 segnali rilevanti/settimana identificati	Watch se segnali già noti
Sintesi ricerca di mercato	Bassa	A	€0,05 – 0,30	3 report	€0,15 – €0,90	≥ 2 insight non ovvi; fonti verificabili	Stop se solo riassunti ovvi
Preparazione OKR / KPI trimestrali	Bassa	B	€0,08 – 0,20	1 sessione	€0,08 – €0,20	OKR SMART completi; approvati senza revisione strutturale	Watch se OKR vaghi

Sezione 4 – Budget mensile: team di 10 persone

FUNZIONE	COSTO MENSILE BASSO	COSTO MENSILE ALTO	NOTE
Marketing & Comunicazione	€4	€12	Ottimizzabile subito su bozze
Commerciale & Vendite	€5	€18	Attenzione a forecast (Livello A)
Operazioni & Amministrazione	€4	€14	Alto volume, basso costo → buon ROI
HR & Formazione	€3	€8	Spesa contenuta, impatto medio-alto
Management & Strategia	€4	€12	Alta variabilità, uso raro ma critico
TOTALE TEAM (10 persone)	€20/mese	€64/mese	Budget controllato se Livello A usato con criterio

Attenzione al salto di scala: questi numeri si riferiscono a uso umano assistito. Se introduci agenti autonomi (AI che agisce senza supervisione su molti task in parallelo), i costi possono moltiplicarsi di 10-50x rapidamente. Il budget per agenti va gestito separatamente.

Sezione 5 – Formula budget settimanale

```
Budget_settimanale_flusso =  
(Volume_task/settimana × Costo_medio_task) + Buffer_20%
```

Esempio – uso assistito:

- Task: sintesi riunioni · Volume: 15/settimana · Livello C: €0,01/task
- Budget: $(15 \times €0,01) \times 1,20 = \mathbf{€0,18/settimana}$ → praticamente gratis

Esempio – agente autonomo (attenzione):

- Agente commerciale: 500 lead/settimana · Livello B · 3 chiamate API per lead
- Costo: $500 \times 3 \times €0,05 = \mathbf{€75/settimana}$ → **€3.900/anno**
- Serve KPI forte per giustificarlo

Sezione 6 — Soglie ROI: Go / Watch / Stop

FORMULA ROI OPERATIVO

$$\text{ROI_AI} = (\text{Valore_tempo_risparmiato} + \text{Valore_errori_evitati} + \text{Valore_output_extra}) / \text{Costo_AI_totale}$$

Come stimare il valore del tempo: $\text{Ore_risparmiate/mese} \times \text{Costo_orario_dipendente}$

Esempio: 10 ore risparmiate \times €35/ora = €350 di valore mensile. Se costo AI del flusso = €50/mese \rightarrow ROI = 350/50 = **7x** \rightarrow Go deciso.

SOGLIE DECISIONALI

ROI OPERATIVO	DECISIONE	AZIONE
$\geq 3,0$	 Go — tieni e scala	Amplia il flusso, aumenta il volume
2,0 – 2,9	 Go — mantieni	Monitora mensile, ottimizza dove possibile
1,2 – 1,9	 Watch — ottimizza	Rivedi stack (abbassa livello?) o aumenta volume
< 1,2	 Stop — refactoring	Cambia approccio: automazione, livello inferiore, o elimina

Soglia minima di sicurezza: ROI \geq 2,0. Sotto quel valore, il costo AI sta erodendo margine senza creare valore netto.

Sezione 7 — Dashboard di controllo (revisione mensile)

30 minuti al mese. Una riga per flusso.

FLUSSO	OWNER	VOLUME/MESE	STACK (LIVELLO)	COSTO/MESE REALE	VALORE STIMATO	ROI	DECISIONE
Sintesi riunioni	Ops	40	C	€0,60	€140	23x	● Go
Bozze newsletter	Mktg	20	B	€1,20	€80	67x	● Go
Analisi forecast	Sales	2	A	€0,60	€200	33x	● Go
Outreach email	Sales	200	B	€6	€120	20x	● Go
[Tuo flusso 1]							
[Tuo flusso 2]							
[Tuo flusso 3]							

Sezione 8 – I 5 errori più comuni

ERRORE	PERCHÉ SUCCEDE	SOLUZIONE
Usare il modello premium su tutto	"Voglio la qualità migliore" → manca fiducia nei livelli inferiori	Test su 20 task: spesso Livello B o C è sufficiente per il 70% dei casi
Nessun KPI definito	Il flusso "sembra utile" ma nessuno misura niente	Prima di attivare un flusso, definire almeno 1 KPI numerico
Budget per sessione, non per flusso	Si paga in modo opaco senza visione d'insieme	Aggregare i costi per flusso, non per singola chiamata
Agenti autonomi senza cap	L'agente scala i task da solo → costi esplodono	Impostare sempre un cap giornaliero e un alert a soglia
Dimenticare il costo del tempo umano	Si conta solo il costo AI, non il tempo di supervisione	Aggiungere il costo della supervisione nella formula ROI

Sezione 9 – Checklist governance AI (mensile)

- Tutti i flussi AI attivi elencati nella dashboard (Sezione 7)
 - Ogni flusso ha un owner e un KPI numerico
 - Costo mensile totale AI calcolato e confrontato con il mese precedente
 - ROI verificato per ogni flusso → applicata soglia Go/Watch/Stop
 - Flussi in Watch: piano di ottimizzazione definito
 - Flussi in Stop: disattivati o riconfigurati
 - Nessun agente autonomo attivo senza cap di costo giornaliero
 - Prossimo mese: un flusso da esplorare, uno da ottimizzare
-

Glossario

TERMINE	SIGNIFICATO PRATICO
Token	Unità di misura del testo elaborato dall'AI. Circa 750 parole = 1.000 token. Si paga in base ai token consumati.
Modello	Il "motore" AI (es. Claude Opus, GPT-4, Gemini). Più è potente, più costa.
Agente AI	Sistema AI che esegue sequenze di azioni in autonomia. Consuma molti più token di una singola chiamata.
Livello A/B/C/D	Classificazione operativa usata in questo documento per evitare nomi tecnici.
KPI minimo	La soglia sotto cui il flusso non sta producendo abbastanza valore. Se non raggiungi il KPI → Watch o Stop.
ROI operativo	Valore prodotto diviso costo AI. Soglia minima consigliata: 2,0.
Cap giornaliero	Limite massimo di spesa sull'agente. Indispensabile per evitare sorprese in bolletta.

Appendice — Prezzi modelli AI (Febbraio 2026)

A.1 — MAPPA MODELLO → LIVELLO OPERATIVO

FORNITORE	MODELLO	LIVELLO	INPUT (\$/1M TOKEN)	OUTPUT (\$/1M TOKEN)	NOTE
Anthropic	Claude Opus 4.5 / 4.6	A — Premium	\$5,00	\$25,00	↓ 67% vs Opus 4.1. Migliore per coding e analisi complesse
Anthropic	Claude Sonnet 4.5 / 4.6	B — Standard	\$3,00	\$15,00	Miglior rapporto qualità/prezzo. Top-10 performance
Anthropic	Claude Haiku 4.5	C — Operativo	\$1,00	\$5,00	Veloce, economico, ideale per volume
OpenAI	GPT-4o	B — Standard	\$2,50	\$10,00	Ampia integrazione Microsoft/Azure
OpenAI	GPT-4o-mini	C — Operativo	\$0,15	\$0,60	16× più economico di GPT-4o
OpenAI	o3 (reasoning)	A — Premium	\$2,00	\$8,00	↓ 80% vs lancio. Per ragionamento matematico/logico
OpenAI	GPT-5.2	A+ — Ultra premium	\$75,00	\$300,00	Solo per casi ad altissimo valore economico
Google	Gemini 2.5 Flash	C — Operativo	\$0,30	\$2,50	Context window 1M token
Google	Gemini 2.5 Pro	B — Standard	\$1,25	\$10,00	Forte su multimodale
Google	Gemini 3 Pro Preview	A — Premium	\$2,00	\$12,00	Vince su helpfulness nei test umani
Mistral AI	Mistral Small	C — Operativo	\$0,10	\$0,30	Ottima alternativa economica europea
Mistral AI	Mistral Medium 3	B — Standard	\$0,40	\$2,00	Performance GPT-4 a costo Haiku
Mistral AI	Mistral Large	A — Premium	\$2,00	\$6,00	Alternativa europea a GPT-4o
Mistral AI	Mistral Nemo	D — Automazione	\$0,02	~\$0,05	Costo minimo assoluto via API
DeepSeek	DeepSeek R1	B/A — Ragionamento	\$0,55	~\$2,19	Budget-friendly. Valutare implicazioni data privacy

Self-hosted	Ollama / LM Studio	D – Automazione	~€0	~€0	Costo fisso infrastruttura, zero per chiamata
--------------------	--------------------	-----------------	-----	-----	---

A.2 – COSTO PER TASK: CALCOLO RAPIDO

Stima per task medio (1.500 token input + 500 token output):

LIVELLO	MODELLO	COSTO/TASK SINGOLO	COSTO/TASK AGENTICO (5 CHIAMATE)	COSTO 100 TASK/MESE
A	Claude Opus 4.5	~€0,019	~€0,095	~€9,50
A	GPT-4o	~€0,012	~€0,060	~€6,00
B	Claude Sonnet 4.5	~€0,012	~€0,060	~€6,00
B	Mistral Medium 3	~€0,002	~€0,010	~€1,00
C	Claude Haiku 4.5	~€0,004	~€0,020	~€2,00
C	Gemini 2.5 Flash	~€0,002	~€0,010	~€1,00
C	GPT-4o-mini	~€0,0004	~€0,002	~€0,20
D	Mistral Nemo	~€0,00004	~€0,0002	~€0,02

Conversione: 1 USD ≈ 0,93 EUR (febbraio 2026)

A.3 – LEVE DI OTTIMIZZAZIONE COSTO

LEVA	RISPARMIO	COME FUNZIONA	QUANDO USARLA
Batch API	~50%	Elaborazione asincrona (risultato in max 24h)	Report notturni, classificazioni non urgenti
Prompt Caching	fino a 90%	Riusa il contesto già inviato senza ripagarci	Agenti con system prompt fisso e lungo
Livello inferiore	50-95%	Scendi da A a B, o da B a C dove la qualità è sufficiente	Testa su 20 task: se accettabile, abbassa
Context management	20-40%	Riduci il contesto passato all'AI	Agenti a lunga esecuzione
Routing intelligente	30-60%	Smista il task al livello corretto in automatico	Stack maturi con volume >500 task/giorno

Combinando batch + caching + routing: risparmio possibile del 75-90% rispetto all'uso non ottimizzato.

A.4 – BENCHMARK QUALITÀ (FEBBRAIO 2026)

TASK	MIGLIORE MODELLO	NOTA
Coding / Software engineering	Claude Opus 4.5/4.6	80.9% su SWE-bench
Ragionamento matematico / logico	OpenAI o3	100% su AIME 2025 mathematics
Elaborazione multimodale	Google Gemini 3 Pro	Integrazione nativa Google Workspace
Preferenza umana generale	Google Gemini 3 Pro	Vince su helpfulness nei test blind
Rapporto qualità/prezzo B2B	Claude Sonnet 4.5 / Mistral Medium 3	Performance top-10 a costo B
Volume massimo (budget basso)	GPT-4o-mini / Gemini 2.5 Flash / Mistral Small	Quasi equivalenti su task semplici

Regola operativa 2026: non esiste un modello "migliore in assoluto". Il 70-80% dei task aziendali può essere gestito da modelli di Livello B o C con risultati identici al Livello A.

AI MASTERY ITALIA

Il Costo Reale dell'AI

GUIDA OPERATIVA PER CHI DECIDE IL BUDGET AI.

FEBBRAIO 2026

Michele Pastorello
michelepastorello.ai