

The logo for Zendata, featuring the word "zendata" in a white, lowercase, sans-serif font. The letter "o" is replaced by a stylized circular icon consisting of three concentric circles. The background is a solid blue color with a subtle pattern of thin, white, intersecting lines that create a grid-like effect.

zendata

Roma, 09/07/2025

Università Campus Bio-Medico di Roma

All'attenzione del Dott. Marco Graziani

Via Álvaro del Portillo, 21

00128 Roma (RM)

Oggetto: Proposta progettuale per piattaforma AI e assistente virtuale a supporto del personale UCBM

Sommario

Chi siamo	4
1. Introduzione al progetto	5
1.1 Panoramica	5
1.2 Obiettivi Dashboard	5
1. Gestione dei Contenuti	5
2. Analisi e Monitoraggio	5
3. Sicurezza e Controllo Accessi	5
1.3 Obiettivi Assistente Virtuale	6
1. Realizzazione di un RAG System avanzato	6
2. Sicurezza by Design	6
2. Architettura dell'assistente virtuale	7
2.1 Panoramica dell'Architettura	7
2.2 Autenticazione e Identità	8
2.3 Layer di Ingestion	8
2.4 Layer di storage	9
2.5 Layer di AI	10
2.6 Layer Applicativo	10
3. Architettura della Dashboard	12
3.1 Panoramica dell'Architettura	12
1. Frontend Amministrativo	12
2. Backend API	12
3. Function App per l'elaborazione documenti	12
4. Storage condiviso	12
3.2 Autenticazione e Controllo Accessi	13
3.3 Infrastruttura di Rete	13
3.4 Raccolta e Analisi delle Metriche	13
4. Struttura dell'offerta	14
4.1 Simulazione dei costi AI	15
4.2 Scenario Self-hosted	16
4.3 Condizioni di pagamento	16
4.4 Validità, licenza e condizioni generali	16

Chi siamo

ZenData è una startup dinamica e innovativa, specializzata nello sviluppo di soluzioni basate sull'intelligenza artificiale, per aumentare l'efficienza e stimolare la crescita

Le nostre soluzioni abbracciano una vasta gamma di settori e reparti, spaziando dallo sviluppo di assistenti virtuali per pubbliche amministrazioni e professionisti autonomi, fino all'analisi e all'ottimizzazione dei processi produttivi aziendali.

La nostra strategia è orientata alla progettazione di soluzioni che si integrano in maniera semplice e fluida nei processi operativi dei nostri clienti, garantendo un'adozione senza attriti e un impatto positivo immediato.

Siamo convinti che il nostro approccio all'innovazione, unito alla capacità di personalizzare le soluzioni in base alle specifiche esigenze dei nostri clienti, ci permetta di offrire un valore aggiunto significativo e di contribuire al successo dei nostri partner commerciali.



1. Introduzione al progetto

1.1 Panoramica

Il presente progetto prevede la realizzazione di un assistente virtuale dedicato ai dipendenti (Personale Tecnico Amministrativo e Accademia) dell'Università Campus Bio-Medico di Roma, su un'architettura progettata per garantire scalabilità, efficienza e sicurezza nella gestione delle informazioni.

L'assistente virtuale pensato offre supporto nella consultazione della documentazione di natura tecnico amministrativa già caricata sulla intranet UCBM come linee guida di processo, regolamenti e istruzioni operative.

Sarà fruibile sia sulla pagina internet UCBM, previa autenticazione con Identity Provider con protocollo SAML 2.0 o OAuth 2.0, che sulla intranet. Anche quando acceduto presso il sito intranet, l'assistente virtuale dovrà poter essere utilizzato previo login integrato con IDP, almeno sino a quando i due enti Università e Fondazione Policlinico Campus Bio-Medico di Roma condivideranno il sito intranet.

A supporto di queste funzionalità, il progetto prevede anche lo sviluppo di una **dashboard**, che consentirà agli amministratori di monitorare l'utilizzo dell'assistente AI, analizzare metriche di performance e di costo, e gestire i contenuti in modo semplice e intuitivo.

L'intera soluzione sarà implementata su tecnologie cloud di Azure, assicurando scalabilità e sicurezza.

1.2 Obiettivi Dashboard

1. Gestione dei Contenuti

- Creazione di un'interfaccia intuitiva per il **caricamento e l'eliminazione dei documenti** della Knowledge Base (KB).

2. Analisi e Monitoraggio

- **Visualizzazione delle metriche di utilizzo**, fornendo insight sulle interazioni degli utenti e sull'efficacia del sistema.
- **Monitoraggio delle performance del sistema RAG**, con l'obiettivo di ottimizzare il recupero delle informazioni e migliorare la qualità delle risposte generate.

3. Sicurezza e Controllo Accessi

- **Integrazione di autenticazione sicura tramite Azure AD** per garantire che solo

utenti autorizzati possano accedere alle funzionalità amministrative.

Tenant	Dominio	Stakeholder
unicampus-int. it	@unicampus .it	Dipendenti (personale tecnico amministrativo e accademia)

1.3 Obiettivi Assistente Virtuale

1. Realizzazione di un RAG System avanzato

- Utilizzo di un **database vettoriale** per la ricerca semantica, migliorando la qualità delle risposte fornite dal chatbot.
- Mantenimento di un **database NoSQL scalabile** basato su **Azure Cosmos DB for MongoDB**, garantendo elevate prestazioni e affidabilità.
- Esecuzione della logica applicativa in un **servizio gestito** tramite **Azure App Service**, assicurando scalabilità e semplificazione della gestione infrastrutturale.

2. Sicurezza by Design

- **Gestione centralizzata** delle identità e degli accessi tramite **Azure AD**, assicurando un controllo rigoroso sugli utenti abilitati.
- **Limitazione del traffico** di rete attraverso **Virtual Network e private endpoints**, proteggendo le comunicazioni interne al sistema.
- **Protezione di segreti e chiavi** sensibili con **Azure Key Vault**, riducendo i rischi legati alla gestione delle credenziali. Le chiavi saranno gestite con versioning abilitato e audit log completo.
- Gestione delle vulnerabilità e patching del sistema, attraverso un piano di aggiornamenti periodici e/o straordinari, in particolare in ordine alle vulnerabilità di sicurezza di livello critico (CVSS score 9.0 – 10.0) che dovessero emergere in qualsiasi dei componenti utilizzati dall'applicativo.
- **Monitoraggio continuo** delle performance e della sicurezza tramite **Azure Monitor**, garantendo proattività nella gestione operativa.

2. Architettura dell'assistente virtuale

La seguente sezione descrive l'architettura prevista per la messa in produzione dell'assistente virtuale. L'infrastruttura sarà implementata su Microsoft Azure, garantendo scalabilità, sicurezza e flessibilità. L'adozione di una piattaforma cloud permetterà di ottimizzare le prestazioni, semplificare la gestione e garantire un'evoluzione continua della soluzione senza impatti sull'operatività.

2.1 Panoramica dell'Architettura

L'architettura si compone sostanzialmente di 4 layer funzionali:

1. **Layer di Ingestion:** rappresenta le soluzioni software e le risorse cloud che permetteranno l'ingestion semi-automatica dei documenti come knowledge base del sistema.
2. **Layer di Storage:** si compone dei database a supporto del funzionamento dell'assistente virtuale. Per questo livello verranno utilizzati un database vettoriale per la ricerca semantica tra la documentazione (Retrieval del RAG system) e un database NoSQL (Azure Cosmos DB for MongoDB) che semplifica la gestione operativa delle conversazioni e al contempo la gestione di collection dedicate alle analytics della dashboard.
3. **Layer di AI:** Questo livello si basa sui Foundational Models forniti da Azure, utilizzati sia per la ricerca documentale tramite modelli di embedding sia per la generazione delle risposte attraverso Large Language Models. I modelli impiegati saranno quelli del servizio Azure OpenAI, di cui di seguito il trattamento dati: [Data, privacy, and security for Azure OpenAI Service](#).
4. **Layer Applicativo:** Si compone delle risorse cloud per la realizzazione delle App backend e frontend.

2.2 Autenticazione e Identità

Per garantire la protezione dei dati e il controllo degli accessi, il sistema adotterà **Azure Active Directory (Azure AD)** per la gestione centralizzata delle identità e l'autenticazione sicura alle risorse Azure.

L'uso delle **Managed Identities** permette ai servizi Azure di autenticarsi in modo sicuro senza necessità di gestire manualmente credenziali statiche. In particolare, l'**App Service** utilizzerà una **System-Assigned Managed Identity** per ottenere token da Azure AD e accedere in modo sicuro a **Key Vault**, **Cosmos DB** e al **database vettoriale**, sfruttando il **Role-Based Access Control (RBAC)** e politiche di accesso dedicate.

Per rafforzare la sicurezza della rete e ridurre la superficie di attacco, il sistema sarà isolato all'interno di una **Virtual Network (VNet) dedicata**. Questa configurazione permette di limitare l'accesso ai soli sistemi autorizzati e instradare il traffico attraverso connessioni private, evitando esposizioni pubbliche. Sarà prevista l'integrazione con la stessa Vnet per l'assistente virtuale e la dashboard.

L'architettura prevede:

- **VNet dedicata**
- **Integrazione dell'App Service con la VNet**, per connettersi internamente a **Cosmos DB**, il **database vettoriale** e **Key Vault** tramite IP privati.
- **Private endpoint per Cosmos DB e il database vettoriale**, assegnati a una subnet dedicata per garantire che il traffico resti confinato all'interno della VNet, senza esposizione pubblica.
- **Private endpoint per Key Vault**, che limita l'accesso solo ai servizi interni alla rete privata.

Questa strategia garantisce un'elevata sicurezza dell'infrastruttura, proteggendo i dati e minimizzando i rischi legati agli accessi non autorizzati.

2.3 Layer di Ingestion

Il layer di ingestion si compone di tutte quelle risorse che hanno l'obiettivo di permettere il popolamento della Knowledge base dell'assistente virtuale.

Questo layer prevede l'utilizzo di Azure Functions per l'estrazione del contenuto della documentazione da integrare all'interno dell'assistente AI. In particolare, queste funzioni si occuperanno di:

- Manipolare il testo non strutturato della documentazione:
 - Preprocessing per l'estrazione del testo e dei metadati utili all'intelligenza artificiale per l'ottimizzazione delle ricerche.
 - Suddivisione in chunk per ottimizzare le prestazioni della RAG.
 - Calcolo dei vettori di embedding per l'indicizzazione dei documenti all'interno del database vettoriale.
- Caricare i documenti pre processati all'interno del database vettoriale.
- Caricare i documenti su un Blob Storage.

2.4 Layer di storage

Il layer di storage consiste sostanzialmente di due database:

- **Azure Cosmos DB for MongoDB** viene utilizzato come database scalabile per la gestione dei dati transazionali utilizzati dall'assistente per il recupero dei messaggi di una conversazione. Contemporaneamente, questa tipologia di Database, permette l'esecuzione di query volte a mostrare le analitiche di utilizzo all'interno di una dashboard di amministrazione.
- Azure AI Search è il database vettoriale adottato per la gestione della documentazione non strutturata alla base del processo di **Information Retrieval (IR)** del RAG system.

Entrambi i database permettono l'accesso in sicurezza ai dati attraverso connessioni su rete privata utilizzando integrazione tramite Virtual Network e l'utilizzo di Private Endpoint.

2.5 Layer di AI

Questo layer si compone dei foundational models che verranno utilizzati sia per la ricerca di documentazione all'interno del database vettoriale (attraverso modelli di Embedding), che per la formulazione delle risposte dell'assistente virtuale in linguaggio naturale (attraverso l'uso di Large Language Models).

Lo stack tecnologico di questo layer prevede:

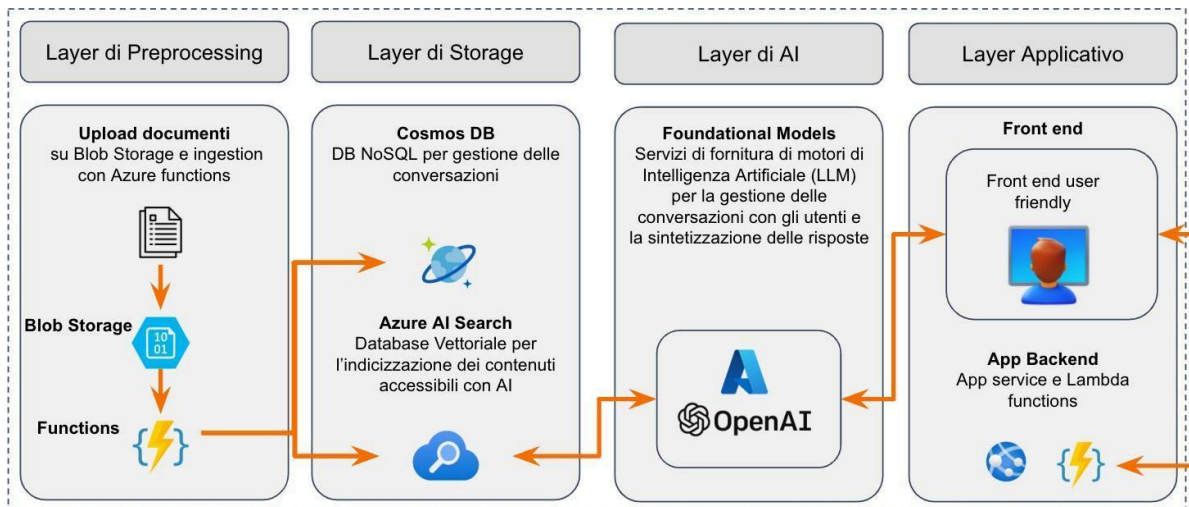
- **Motore di Embedding (IR)** che utilizza la domanda dell'utente per cercare le sezioni della documentazione più pertinenti all'interno del Database vettoriale.
- **Large Language Model** che, utilizzando i contenuti estratti dal motore IR, contestualizza la risposta fornita.

2.6 Layer Applicativo

Per l'hosting e la gestione della logica applicativa, verrà utilizzato **Azure App Service**, che eseguirà l'intero sistema all'interno di **container Docker**. Questa soluzione garantisce scalabilità, isolamento e una gestione semplificata delle risorse cloud. Per il deployment dell'interfaccia front-end è stata individuata la risorsa **Azure Static Web App**. Questa soluzione offre un'elevata integrazione con gli altri servizi Azure e facilita l'inserimento nei siti esistenti. Tuttavia, nel corso del progetto verranno valutati ulteriori dettagli per garantire la massima compatibilità con i sistemi di destinazione. L'accesso del backend applicativo agli altri layer dell'architettura avverrà tramite **Virtual Networks (VNet)** e **Private Endpoints**, garantendo connessioni sicure su rete privata. Le chiavi di accesso saranno archiviate in **Azure Key Vault** e gestite attraverso **Managed Identity**, assicurando un'integrazione sicura e senza l'uso di credenziali esplicite.

Per il monitoring dello stato dell'applicazione verrà adottato Azure Application Insights.

Di seguito si riporta l'architettura ad alto livello precedentemente descritta:



Come requisito modulare sarà possibile integrare interrogazione e risposta in modalità vocale. In particolare utilizzando:

- **Speech-to-text:** utilizzo di Azure Speech Service o altri servizi per la trascrizione vocale in tempo reale
- **Text-to-speech:** utilizzo di Azure Speech Service o altri servizi per la sintesi vocale

I dati vocali saranno archiviati temporaneamente per il tempo strettamente necessario alla transcodifica, e successivamente cancellati in modo sicuro, con policy documentabili.

3. Architettura della Dashboard

La seguente sezione descrive l'architettura prevista per la messa in produzione della Dashboard. L'infrastruttura sarà implementata su Microsoft Azure, garantendo perfetta integrazione con l'assistente virtuale precedentemente descritto.

L'architettura della Dashboard è stata progettata per offrire un'interfaccia amministrativa completa che consenta la gestione efficiente dei contenuti e il monitoraggio delle performance dell'assistente virtuale.

3.1 Panoramica dell'Architettura

L'architettura della Dashboard si compone di quattro componenti principali:

1. Frontend Amministrativo

- Offre funzionalità di gestione dei documenti e visualizzazione delle analytics
- Design responsive ottimizzato per diversi dispositivi

2. Backend API

- Gestisce le richieste dal frontend e coordina le operazioni sui dati

3. Function App per l'elaborazione documenti

- Componente serverless che gestisce operazioni asincrone e di lunga durata
- Si occupa dell'estrazione del testo dai documenti e della creazione di chunk ottimizzati

4. Storage condiviso

- Utilizzo delle stesse risorse di storage dell'assistente virtuale
- Include Cosmos DB, database vettoriale e Blob Storage
- Garantisce coerenza dei dati e riduce la duplicazione

Le attività principali supportate da questa architettura includono:

- Gestione completa del ciclo di vita dei documenti (caricamento, elaborazione, indicizzazione, eliminazione)
- Monitoraggio delle performance dell'assistente virtuale attraverso metriche e analytics
- Amministrazione sicura con autenticazione Azure AD. Le attività amministrative saranno tracciate tramite audit trail persistente (immutabile), consultabile dal personale di Ateneo.

3.2 Autenticazione e Controllo Accessi

Per garantire la protezione dei dati e il controllo degli accessi, il sistema adotterà le stesse politiche di sicurezza descritte nell'architettura dell'assistente virtuale. Azure Active Directory (Azure AD) viene utilizzato per autenticare gli amministratori e i servizi che accedono alla dashboard, garantendo un accesso centralizzato e sicuro. L'uso delle Managed Identities consente all'applicazione di autenticarsi automaticamente con i servizi Azure senza necessità di gestire manualmente credenziali.

3.3 Infrastruttura di Rete

Verrà utilizzata una subnet dedicata per i componenti della Dashboard all'interno della VNet esistente. Verranno utilizzati gli stessi private endpoint configurati per l'assistente virtuale, assicurando che il traffico resti confinato all'interno della rete privata.

3.4 Raccolta e Analisi delle Metriche

La dashboard prevede la raccolta di metriche pensate per monitorare le performance, individuare trend e ottimizzare il servizio nel tempo. Le metriche definitive potranno essere accordate in dettaglio in una seconda fase, ma includeranno il tracciamento delle query al chatbot (volume, tipologia, distribuzione temporale) e l'analisi delle tematiche più frequenti. I dati raccolti verranno poi resi disponibili attraverso dashboard interattive.

4. Struttura dell'offerta

Di seguito si riporta il dettaglio dei costi complessivi di progetto che comprende la realizzazione e la manutenzione dell'assistente virtuale e della dashboard amministrativa:

Costi di Progetto	Totale	Sconto	Totale scontato
implementazione (CapEx)	10.000 €	20%	8.000 €
manutenzione (OpEx) (esclusa infrastruttura)	500 €/mese	14%	430 €/mese

Di seguito si riportano i costi relativi all'infrastruttura cloud necessaria per l'erogazione della soluzione, validi per entrambe le modalità di fornitura (SaaS o self-hosted su Azure):

Nome Servizio	Descrizione	€/mese
Azure Cosmos DB	Azure Cosmos DB for MongoDB (RU), Autoscale provisioned throughput	40
Azure AI Search	Basic, 1 Unit(s), 730 hours	70
Virtual Network	Sweden Central region	10
Azure Blob Storage	Block Blob Storage, Standard, General Purpose V2, Flat Namespace, Hot, LRS, 50 GB,	2

Azure functions	Sweden Central region, Premium	150
App Service	Premium V3, 2 core, 8 GB RAM, 250 GB Storage. Backend dell'assistente virtuale. 1 year reserved	90
App Service	Basic B2, 3.5 GB RAM 10GB Storage. Backend Dashboard amministrativa. On demand.	20
Key Vault	Sweden Central region	1
Totale		383
Sconto		14%
Totale Scontato		329,38

4.1 Simulazione dei costi AI

Oltre ai costi infrastrutturali fissi sopra riportati, l'assistente virtuale comporta costi operativi legati all'utilizzo dell'intelligenza artificiale, che variano in base all'effettivo volume di utilizzo.

A supporto dell'offerta, è incluso il documento Excel denominato

calcolatore_costi_chatbot_UCBM.xlsx, contenente un foglio di calcolo che consente di simulare in modo dettagliato questi costi. In particolare il file permette di stimare i costi mensili in funzione di vari parametri personalizzabili, tra cui:

- numero di utenti;
- numero medio di conversazioni mensili;
- numero di messaggi per conversazione;
- percentuale di interazioni vocali.

Tutti i costi sono espressi su base mensile e possono essere adattati in base a diversi scenari di utilizzo, consentendo così una proiezione trasparente e flessibile nel tempo.

Il file può essere utilizzato per analizzare la sostenibilità economica del sistema e confrontare facilmente i costi in diverse configurazioni operative.

4.2 Scenario Self-hosted

I costi infrastrutturali indicati rimangono invariati sia nel caso di fornitura in modalità SaaS, sia qualora l'infrastruttura venga ospitata direttamente da UCBM (scenario self-hosted), a condizione che venga utilizzata la piattaforma cloud Microsoft Azure.

Nel caso si optasse per la modalità self-hosted, Zendata si impegna a fornire la documentazione tecnica e un handover completo del codice e delle configurazioni.

4.3 Condizioni di pagamento

- **CAPEX:** 90 gg DF.
- **OPEX:** Semestrale.

4.4 Validità, licenza e condizioni generali

- **Validità dell'offerta:** 60 giorni dalla data di emissione.
- **Durata della licenza d'uso:** 12 mesi dalla data di attivazione del sistema.
- **Modalità di rinnovo:** Rinnovo annuale opzionale alle stesse condizioni, salvo modifiche concordate.
- **Tempistiche di consegna:** entro 6 settimane dalla conferma dell'ordine e sottoscrizione contrattuale.

Licenza d'uso e proprietà intellettuale

- La licenza è non esclusiva, non trasferibile, e limitata all'uso da parte di UCBM per gli scopi definiti nel presente progetto.
- Tutti i diritti sul codice sorgente, le architetture e le configurazioni restano di proprietà di Zendata.
- In modalità self-hosted, l'accesso di Zendata all'infrastruttura sarà regolato tramite permessi controllati su Azure.

zendata



info@zendata.it
www.zendata.it

zendata