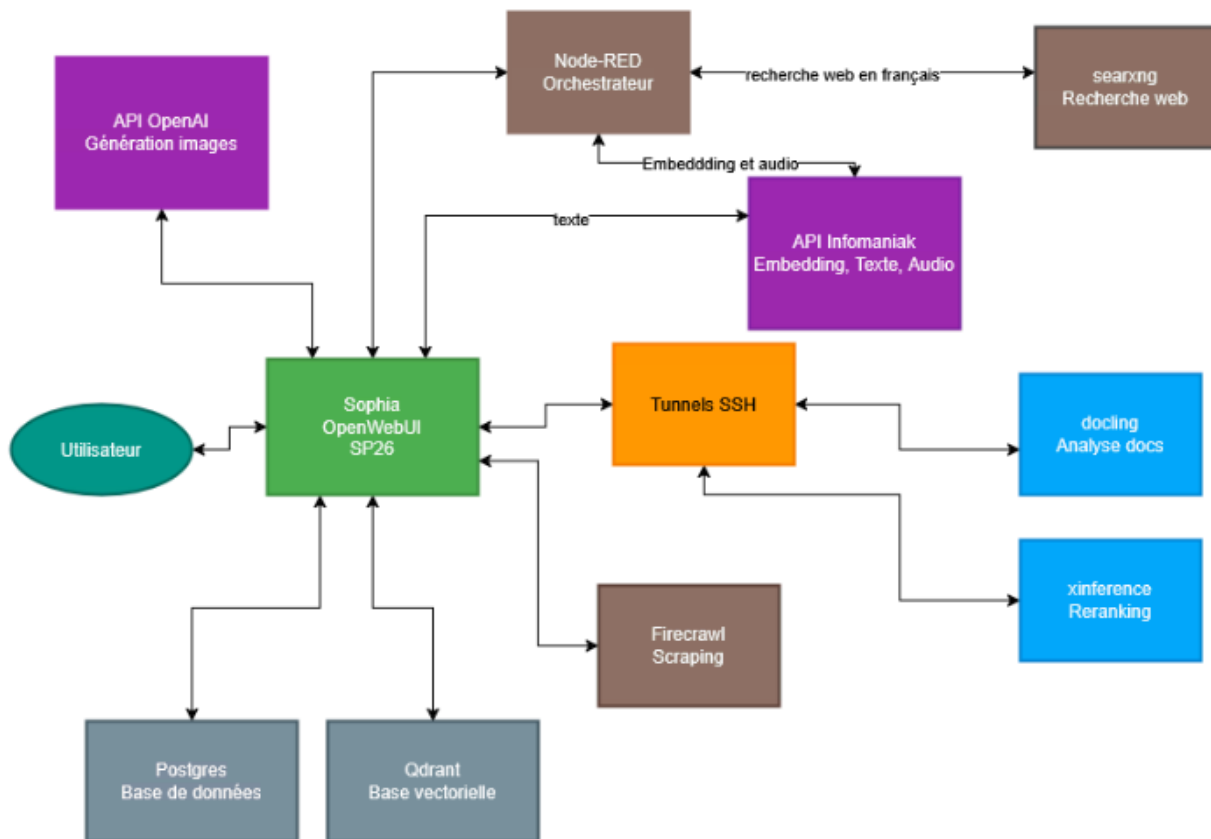


Sophia

- Architecture générale
- Vérification du bon fonctionnement
- Redémarrage de la VM GPU sur Infomaniak

Architecture générale

Schéma d'architecture générale de Sophia



Pour modifier l'image : <https://onlyoffice.vincennes.fr/Products/Files/DocEditor.aspx?fileid=40912>

L'architecture Today at 19:45

L'architecture est basée sur la documentation d'openwebui : <https://docs.openwebui.com/>.

Tableau récapitulatif des conteneurs - visualisation par Portainer

Conteneur / Service	Stack / Origine	Rôle	Port hôte (local)	Port conteneur	URL d'accès depuis Sophia
<code>openwebui</code>	Sophia	Interface principale	3002	8080	<code>https://sophia.vincennes.fr</code>
<code>qdrant</code>	Sophia	Base vectorielle (RAG)	6333, 6334	6333, 6334	Interne (via API locale)
<code>postgres</code>	Sophia	Base de données d'application	5437	5432	Interne
<code>sophia-nodered-openwebui-1</code>	Sophia	Automatisation, intégration des services	1886	1880	<code>http://sp26.vincennes.fr:1886</code>
<code>infomaniak-xinference-1</code>	Infomaniak (tunnel SSH)	Tunnel vers xinference (reranking)	9995	9998	<code>http://sp26.vincennes.fr:9995</code>
<code>infomaniak-docling-1</code>	Infomaniak (tunnel SSH)	Tunnel vers docling (analyse de documents)	5002	5001	<code>http://sp26.vincennes.fr:5002</code>
<code>infomaniak-portainer-agent-1</code>	Infomaniak (tunnel SSH)	Agent Portainer pour supervision distante	9002	9001	Interne (gestion Docker)
<code>searxng</code>	ia-app	Moteur de recherche web	8056	8080	Via Node-RED (recherche)
<code>firecrawl-api-1</code>	firecrawl	Analyse de pages web (scraping intelligent)	3003	3002	<code>http://sp26.vincennes.fr:3003</code>
<code>docling</code> (VM Infomaniak)	-	Analyse de documents (PDF, etc.)	5001	5001	Accès distant (via tunnel SSH depuis SP26:5002)
<code>xinference</code> (VM Infomaniak)	-	Modèle de reranking (<code>bge-reranker-v2-m3</code>)	9997	9997	Accès distant (via tunnel SSH depuis SP26:9995)
<code>portainer-agent</code> (VM Infomaniak)	-	Agent Docker pour supervision	9001	9001	Accès distant (via tunnel SSH depuis SP26:9002)

Réglages dans l'interface administrateur - visualisation par Openwebui

1. Menu Connexions

- **URL** : `https://sophia.vincennes.fr/admin/settings/connections`
- **Fonction** : Gestion des connexions aux modèles IA hébergés sur Infomaniak.
- **Détails** :
 - Affiche l'état de disponibilité des endpoints IA (texte, audio, image).
 - Aucune action de configuration possible directement : en cas de problème, contact avec le support requis.
 - Vérification de la connectivité aux API distantes.

2. Menu Documents

- **URL** : `https://sophia.vincennes.fr/admin/settings/documents`
- **Fonction** : Configuration du traitement des documents (analyse, embedding, reranking).
- **Paramètres** :
 - **Analyse de documents** :
 - URL : `http://sp26.vincennes.fr:5002/`
 - Connecté au service **docling** via tunnel SSH.
 - **Embedding** :
 - URL : `http://sp26.vincennes.fr:1886/`
 - Passerelle via **Node-RED**, qui relaye vers le modèle `bge_multilingual_gemma2` sur Infomaniak.
 - Utilisation de Node-RED pour limiter et réguler les appels API.
 - **Reranking (RAG)** :
 - URL : `http://sp26.vincennes.fr:9995/v1/rerank`
 - Connecté au modèle `bge-reranker-v2-m3` via **xinference** sur Infomaniak (tunnel SSH).

3. Menu Recherche web

- **URL** : `https://sophia.vincennes.fr/admin/settings/web`
- **Fonction** : Configuration de la recherche web et de l'analyse des résultats.
- **Paramètres** :
 - **Recherche** :
 - URL : `http://sp26.vincennes.fr:1886/search?q=<query>`
 - Passerelle via **Node-RED**, qui interroge **searxng**.
 - Node-RED assure la traduction de la requête en français si nécessaire.
 - **Analyse des pages** :
 - URL : `http://sp26.vincennes.fr:3003/`
 - Envoi des URLs récupérées vers **Firecrawl** pour extraction de contenu structuré.

4. Menu Interface utilisateur

- **Fonction** : Génération automatique de métadonnées pour améliorer l'expérience utilisateur.
- **Utilisation** :
 - Appels aux **API Infomaniak** pour :
 - Génération de **titres** pertinents.
 - Création de **tags** et de **mots-clés** pour la recherche.
 - Ces fonctionnalités enrichissent les réponses et les contenus générés.

5. Menu Audio

- **URL** : `https://sophia.vincennes.fr/admin/settings/audio`
- **Fonction** : Transcription de fichiers audio en texte.
- **Paramètres** :
 - **Transcription** :
 - URL : `http://sp26.vincennes.fr:1886/`
 - Utilisation de **Node-RED** pour gérer l'appel à l'API Infomaniak.
 - Le processus est asynchrone : Node-RED surveille la tâche et renvoie le texte une fois disponible.
 - Permet de traiter des fichiers audio longs via une file d'attente.

6. Menu Images

- **Fonction** : Génération d'images à partir de descriptions textuelles.
- **Paramètres** :
 - Appel direct aux **endpoints OpenAI** (ex: DALL·E).
 - Aucun traitement local ou intermédiaire : la requête est envoyée directement à OpenAI.
 - Les images générées sont intégrées dans les réponses de l'interface.

Vérification du bon fonctionnement

Voici une **checklist de vérification du bon fonctionnement** du système, suivie d'un **tableau récapitulatif des problèmes possibles en cas d'indisponibilité**.

☐ Checklist de vérification du bon fonctionnement

Élément à vérifier	Procédure de vérification	Statut (✓/X)
1. Containers dans Portainer	Accéder à l'interface Portainer et s'assurer que tous les conteneurs sont en état "Running" : - <code>openwebui</code> , <code>qdrant</code> , <code>postgres</code> , <code>sophia-nodered-openwebui-1</code> - <code>infomaniak-xinference-1</code> , <code>infomaniak-docling-1</code> , <code>infomaniak-portainer_agent-1</code> - <code>searxng</code> , <code>firecrawl-api-1</code> (et les conteneurs de la stack firecrawl)	
2. Interface Sophia (OpenWebUI)	Accéder à : <code>https://sophia.vincennes.fr</code> → Vérifier que la page charge correctement et que l'interface est interactive.	
3. Interface Node-RED	Accéder à : <code>http://sp26.vincennes.fr:1886</code> → Vérifier que l'éditeur s'ouvre et que les flux sont actifs.	
4. Firecrawl (API et message par défaut)	Accéder à : <code>http://sp26.vincennes.fr:3003</code> → Vérifier que la réponse "SCRAPERS-JS: Hello, world! K8s! "s'affiche.	

Élément à vérifier	Procédure de vérification	Statut (✓/X)
5. Interface Docling	Accéder à : <code>http://sp26.vincennes.fr:5002/ui</code> → Vérifier que l'interface utilisateur de Docling est accessible.	
6. XInference et modèle de reranking	1. Accéder à l'interface XInference via le tunnel : <code>http://sp26.vincennes.fr:9995</code> 2. Vérifier que le modèle <code>bge-reranker-v2-m3</code> est chargé et en cours d'exécution dans la section <code>/reranker</code> .	
7. SearXNG	Accéder à : <code>http://sp26.vincennes.fr:8056</code> → Vérifier que l'interface de recherche s'affiche. → Effectuer une recherche test (ex: "test").	
8. Endpoint OpenAI	Depuis l'interface Sophia (Menu Images) : → Lancer une génération d'image. → Vérifier que la requête atteint OpenAI et retourne une image.	
9. Endpoints Infomaniak (via Node-RED)	Vérifier dans les menus Sophia : - Connexions : état des endpoints texte/audio/image - Documents, Audio, Recherche : que les appels passent via <code>http://sp26.vincennes.fr:1886</code> → Vérifier les logs Node-RED en cas d'erreur.	
10. Logs des conteneurs de tunnel SSH	Vérifier les logs des conteneurs suivants dans Portainer : - <code>infomaniak-xinference-1</code> - <code>infomaniak-docling-1</code> - <code>infomaniak-portainer_agent-1</code> → Rechercher des erreurs de connexion SSH, timeout ou refus. → S'assurer que la VM Infomaniak est joignable (ping ou test de connexion réseau).	

⚠ Tableau des problèmes possibles en cas d'indisponibilité

Service indisponible	Symptôme observé	Cause probable	Impact fonctionnel
Conteneur non démarré dans Portainer	Interface inaccessible, erreur 502/503	Conteneur crashé, mauvaise configuration, dépendance manquante	Blocage total du service associé
Sophia (OpenWebUI)	Page blanche, erreur de chargement	Problème réseau, conteneur down, certificat SSL expiré	Interface principale inutilisable
Node-RED	Éditeur inaccessible ou flux inactifs	Conteneur down, erreur dans un flux critique	Toutes les intégrations (recherche, audio, embedding) bloquées
Firecrawl	API ne répond pas ou timeout	Conteneur down, VM surchargée, erreur de scraping	Impossible d'extraire du contenu web pour les recherches
Docling	/ui inaccessible ou erreur 500	Tunnel SSH rompu, service Docling arrêté sur Infomaniak	Analyse des documents PDF/HTML impossible
XInference / bge-reranker-v2-m3	Modèle non chargé ou erreur 404/500	Modèle non démarré, mémoire insuffisante, tunnel SSH défaillant	Détérioration de la qualité des réponses (RAG)
SearXNG	Page de recherche inaccessible	Conteneur down ou configuration réseau incorrecte	Recherche web impossible via Sophia
Endpoint OpenAI	Génération d'image échoue	Clé API invalide, réseau bloqué, quota dépassé	Fonctionnalité d'image inactive
Endpoints Infomaniak	Erreurs dans les menus Sophia (Connexions, Audio, etc.)	Tunnel SSH rompu, VM éteinte, problème d'authentification	Toutes les fonctions IA distantes (transcription, embedding) inactives
Logs des tunnels SSH	Messages répétés : Connection refused, Broken pipe, Failed to connect	VM Infomaniak éteinte, firewall, clé SSH expirée, réseau instable	Perte de contrôle et de données vers les services distants

Recommandation : Cette checklist peut être utilisée quotidiennement ou après chaque déploiement/mise à jour. En cas d'indisponibilité, consulter en priorité les **logs des conteneurs de tunnel SSH** et vérifier la **connectivité réseau vers la VM Infomaniak**.

Redémarrage de la VM GPU sur Infomaniak

Aller sur <https://api.pub2.infomaniak.cloud/horizon/project/>

Les identifiants sont dans psono "Infomaniak openstack / Horizon"

Puis sur la page <https://api.pub2.infomaniak.cloud/horizon/project/instances/>

Une seule instance est configurée : "whisperx".

Instances

ID de l'instance = <input type="text"/>												Filter	Lancer une instance	Supprimer les instances	Plus d'actions ▾						
Affichage de 1 élément																					
<input type="checkbox"/>	Instance Name	Image Name	IP Address	Flavor	Key Pair	Status	Availability Zone	Task	Power State	Age	Actions										
<input type="checkbox"/>	whisperx	-	37.156.45.74, 2001:1600:16:10::209	nvl4-a4-ram8-disk50-perf2	Vincennes	Active	az-1	Aucun	En fonctionnement	6 mois	Créer un instantané ▾										

Dans le cas où la VM présentera un problème de blocage persistant, elle peut être redémarrée via l'option disponible en bout de ligne : "Redémarrer matériellement l'instance".

Le redémarrage prend 2 à 3 minutes.

Une fois relancée, les containers réapparaissent fonctionnels sur Portainer / Infomaniak.

Les accès à docling et xinference doivent être vérifiés (voir points 5 et 6 de <https://formagent.vincennes.fr/books/intelligence-artificielle/page/verification-du-bon-fonctionnement>).