

Bloc de Compétences E1:

Data Platform IMO-Ops

Réalisation de la collecte, le stockage et la mise à disposition des données pour un projet en intelligence artificielle dédié au marché Immobilier.

Rédigé par : Olivier LAVAUD, le 14 Novembre 2025

Pour l'acquisition du Titre RNCP 37827:
Développeur en Intelligence Artificielle

<https://github.com/OlivierLAVAUD/imo-ops.git>

Microsoft en France

Ecole IA

DEVELOPPEUR·SE EN
INTELLIGENCE
ARTIFICIELLE

SIMPLON



Data Platform IMO-Ops pour un projet en intelligence artificielle dédié au marché Immobilier.

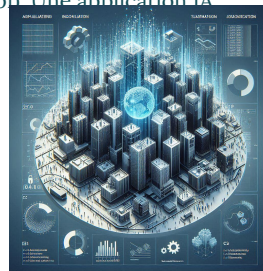
Dev IA RNCP 37827 / © 2025 Olivier LAVAUD

Contexte & objectifs:

- L'objectif est de proposer une **Data Platform** du marché immobilier alimentée par différentes sources (fichiers, web databases, api). Une application IA par exemple un chatbot LLM RAG pourra avoir accès à des données en temps réel structurées et enrichies.

L'architecture (écosystème/stack technologique)

- Orchestration de pipelines avec messaging, parallélisation & mise à l'échelle (Apache Airflow+Redis - Dags - Logs)
- Scraping: Playwright avec fichiers de configuration
- Manipulation et transformation des données: Pandas, Regex, SQLAlchemy, Requests
- Bases de données: Sql (PostgreSQL-pgadmin) & NoSQL (MongoDB)
- sources de fichiers batch: (JSON, ...)
- sources à partir d'APIs: FastAPI
- Rapports & consultation des données (Prometheus- Grafana, Plotly-Dashboards, Gradio)
- Conteneurisation: Docker



Les processus:

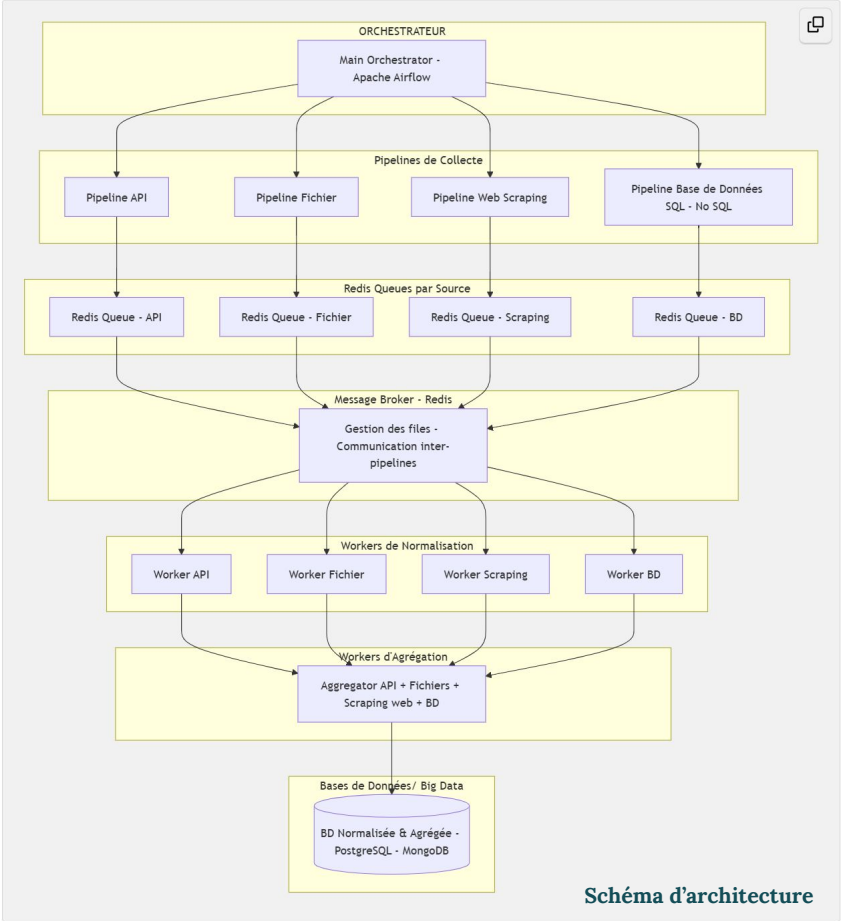
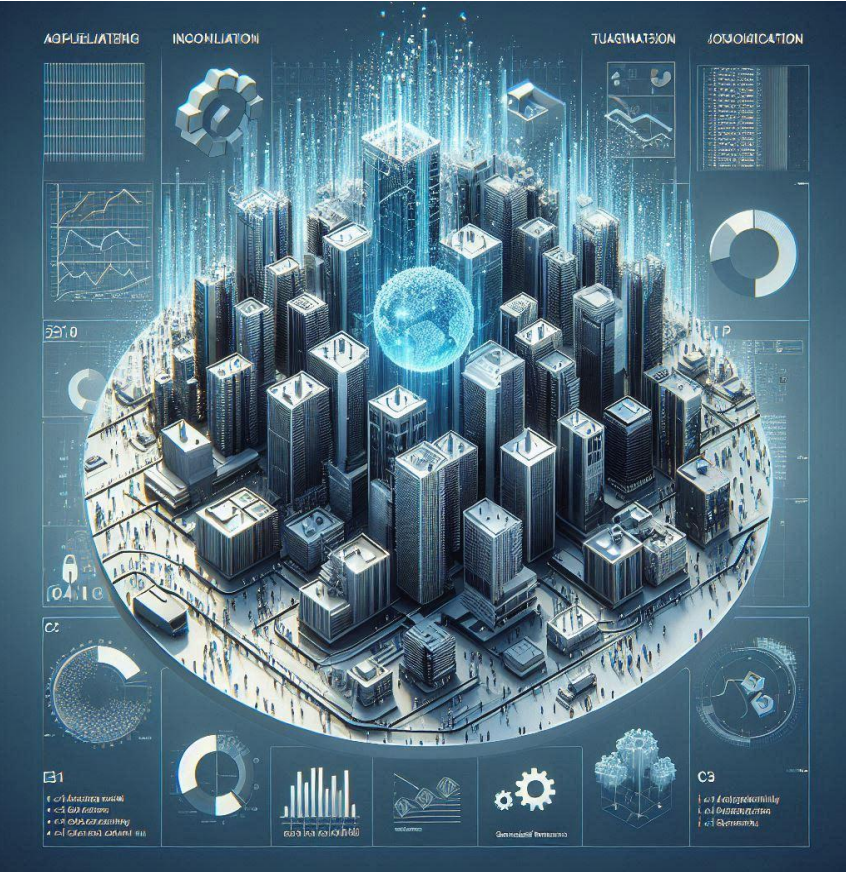
- batches/pipelines de collecte à partir de différentes sources: scraping web, fast API, fichiers csv, bases de données postgres, fichiers csv
- batches/pipelines de normalisation après l'extraction des données
- batch/pipeline d'agrégation des données en provenance des différentes sources
- batch/pipelines de stockage en base de données SQL & NoSQL
- batch/pipeline d'interrogation des données
- batch/pipeline de mise à disposition des données

Les sources:

site web: IAD-Immobilier (<https://www.iadfrance.fr/>): annonces achats, ventes, location sur critères multiples, fichiers, web databases, API

Data Platform IMO-Ops pour un projet en intelligence artificielle dédié au marché Immobilier.

Dev IA RNCP 37827 / © 2025 Olivier LAVAUD



Compétence C1: Automatisation de l'extraction de données

Data from Web Scraping, batch files, sources Databases and APIs

services/c1-scrap (iad-scraper, iad-scraper-custom, extract-from-files, extract-from-api, extract-from-db)

Le scraper web high-tech implémente une architecture asynchrone basée sur Playwright pour le rendu JavaScript complet, avec une gestion avancée du cycle de vie des pages via async/await.

Le système utilise un pattern de stratégies multiples pour le bypass des bannières cookies, combinant sélecteurs CSS dynamiques, injection JavaScript et simulation d'événements claviers.

L'extraction des données suit une approche modulaire avec des classes spécialisées : **DataExtractor** pour le parsing HTML via XPath/CSS combinés, **MediaExtractor** pour l'optimisation des URLs multimédias avec regex de nettoyage qualité.

La configuration externalisée en JSON schema (config.json) permet la définition déclarative des sélecteurs, tandis que le pipeline de pagination intelligent gère automatiquement la navigation multi-pages. L'implémentation inclut un system de retry robuste, des timeouts adaptatifs et la génération de datasets structurés avec métadonnées complètes.

Les autres services de collecte suivent tous la même logique de la normalisation avant agrégation en base de données.

Scénarios courants

```
# Lancement standard (Paris, 200 biens)
docker-compose up iad-scraper
```

```
# Recherche à Toulouse, 25 biens (PowerShell)
$env:LOCALISATION="Toulouse"; $env:MAX_BIENS=25; docker-compose up
iad-scraper-custom
```

```
# Recherche à Strasbourg, 40 biens (Linux/Ubuntu)
LOCALISATION="Strasbourg" MAX_BIENS=40 docker-compose up iad-scraper-custom
```

```
# Mode démon (exécution en arrière-plan)
docker-compose up -d iad-scraper-custom
```

```
# Grandes métropoles
LOCALISATION="Lyon" MAX_BIENS=100 docker-compose up iad-scraper-custom
LOCALISATION="Marseille" MAX_BIENS=80 docker-compose up iad-scraper-custom
```

```
# Villes moyennes
LOCALISATION="Montpellier" MAX_BIENS=50 docker-compose up iad-scraper-custom
LOCALISATION="Nantes" MAX_BIENS=60 docker-compose up iad-scraper-custom
```

```
# Test de développement (limité)
LOCALISATION="Paris" MAX_BIENS=5 MAX_PAGES=1 docker-compose up
iad-scraper-custom
```

```
{
  "site": "iadfrance.fr",
  "name": "IAD France",
  "base_url": "https://www.iadfrance.fr",
  "search_url": "https://www.iadfrance.fr/annonces/vente",
  "wait_timeout": 10000,
  "navigation_delay": 2000,
  "navigation_menu": {
    "container": "nav-menu",
    "css_selector": "nav-menu",
    "xpath": "//nav[contains(class, 'flex-1')]",
    "description": "Menu de navigation principal"
  }
}
```

Compétences C2: Développement de requêtes SQL

SQL Scripts

services/c2-sql (stats_views.sql)

Vues SQL d'Analyse Immobilière Avancée

Les 15 vues SQL constituent un système d'analyse prédéfinie exploitant les fonctionnalités analytiques de PostgreSQL. Chaque vue implémente des agrégations complexes : `v_stats_type_bien` utilise les fonctions d'agrégation standard (AVG, MIN, MAX) avec regroupement multi-critères, tandis que `v_stats_distribution_surfaces` emploie `PERCENTILE_CONT(0.5)` pour calculer les médianes robustes aux valeurs extrêmes.

Les vues temporelles comme `v_stats_prix_m2_par_mois` exploitent `DATE_TRUNC('month')` pour l'analyse de tendances saisonnières, avec partitionnement par type de bien. Les analyses DPE (`v_stats_performance_energetique`) utilisent des jointures optimisées entre les tables annonces et dpe, avec tri sémantique via CASE pour ordonner naturellement les classes énergétiques.

Les vues de distribution comme `v_stats_caracteristiques_populaires` calculent des pourcentages relatifs avec sous-requêtes corrélées, tandis que `v_stats_rapport_surface_prix` implémente des ratios prix/surface avec gestion des divisions par zéro via `NULLIF`. Les catégorisations avancées dans `v_stats_annee_construction` et `v_stats_coproprietes` utilisent des plages dynamiques avec CASE pour segmenter les données en cohorts pertinentes.

La vue synthèse `v_synthese_generale` agrège des métriques cross-tables via sous-requêtes scalaires, offrant un dashboard performance en une seule requête. L'ensemble utilise systématiquement le typage explicite (`::numeric`) et le tri sémantique pour une exploitation analytique immédiate.

Extrait de scripts:

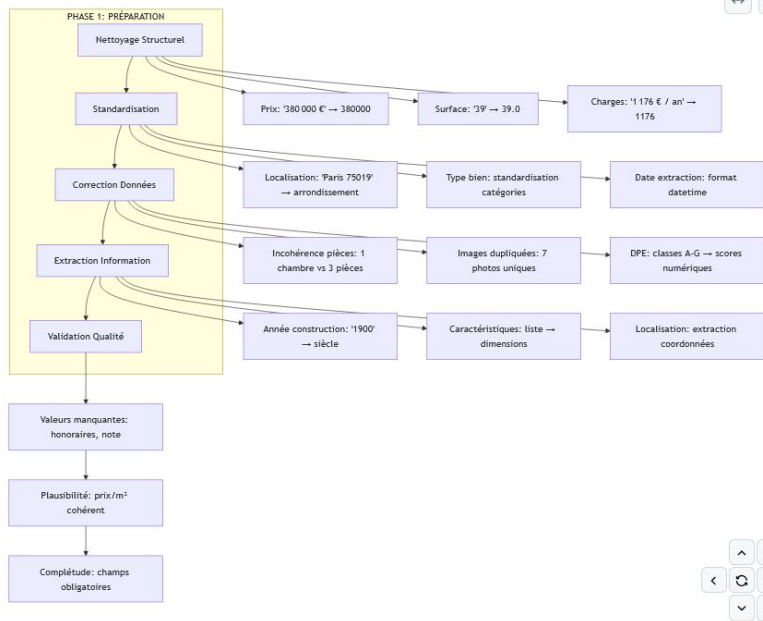
```
-- Segmentation dynamique des prix avec ranking et analyse comparative
WITH ranked_annonces AS (
    SELECT
        localisation,
        type_bien,
        prix,
        prix_au_m2,
        ROUND(AVG(prix_au_m2) OVER (
            PARTITION BY localisation, type_bien
        )::numeric, 2) as prix_m2_moyen_zone,
        ROUND((prix_au_m2 / NULLIF(AVG(prix_au_m2) OVER (
            PARTITION BY localisation, type_bien
        ), 0) - 1) * 100, 2) as ecart_percent_moyen,
        NTILE(4) OVER (
            PARTITION BY localisation, type_bien
            ORDER BY prix_au_m2
        ) as quartile_prix,
        RANK() OVER (
            PARTITION BY localisation
            ORDER BY prix_au_m2 DESC
        ) as rank_prix_dans_commune
    FROM annonces
    WHERE prix_au_m2 IS NOT NULL
      AND localisation IS NOT NULL
      AND type_bien IS NOT NULL
)
SELECT *
FROM ranked_annonces
WHERE rank_prix_dans_commune <= 10 -- Top 10 par commune
ORDER BY localisation, type_bien, prix_au_m2 DESC;
```

	localisation	type_bien	prix	prix_au_m2	prix_m2_moyen_zone	ecart_percent_moyen	quartile_prix	rank_prix_dans_commune
15	Montpellier (34070)	Appartement	340000.00	4000.00	3301.02	21.17	4	5
16	Montpellier (34070)	Appartement	379500.00	3872.45	3301.02	17.31	4	7
17	Montpellier (34070)	Appartement	149000.00	3640.51	3301.02	16.36	4	8
18	Montpellier (34070)	Appartement	199000.00	3764.72	3301.02	13.74	4	10
19	Montpellier (34070)	Duplex	299000.00	3883.12	3883.12	0.00	1	6
20	Montpellier (34070)	Maison	250000.00	4300.00	4076.97	5.47	4	3
21	Montpellier (34070)	Maison	619000.00	4086.61	4076.97	0.24	2	3
22	Montpellier (34070)	Maison	519000.00	4086.61	4076.97	0.24	3	3

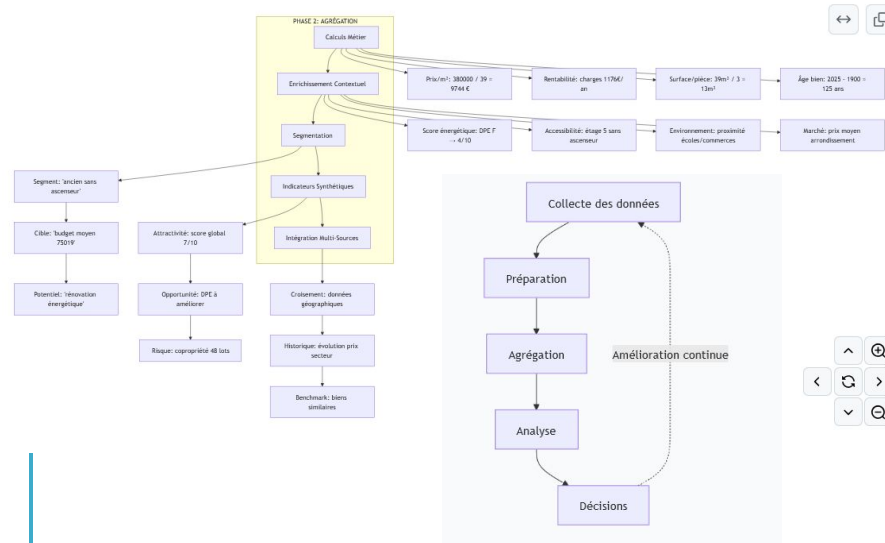
Compétences C3: Agrégation et préparation des données

Data Preparation (Data Cleaning, Data Integration, Data Transformation, Data Reduction) & Data Aggrégation
services/c3_aggr ([README.md](#), README.preparation, README.agregation)

Préparation des données (Data Préparation)



Agrégation des données (Data Aggrégation)



La data preparation est le processus crucial de nettoyage, de transformation et d'enrichissement des données brutes pour les rendre exploitables pour l'analyse.

L'agrégation de données est le processus qui consiste à collecter et rassembler des informations dispersées pour les synthétiser en ensembles cohérents, permettant une analyse globale et une prise de décision éclairée

Compétence C4 : Création de Base de Données

Create and Ingestion in Database

services/c4-create_db

J'ai conçu et implémenté une base de données PostgreSQL complète pour la gestion d'annonces immobilières, en partant d'un modèle UML jusqu'au déploiement opérationnel. Le schéma relationnel intègre six tables principales (annonces, caractéristiques, images, conseiller, DPE, copropriété) avec leurs relations, contraintes d'intégrité et index optimisés. Des vues matérialisées permettent des analyses statistiques sur les prix, surfaces et performances énergétiques.

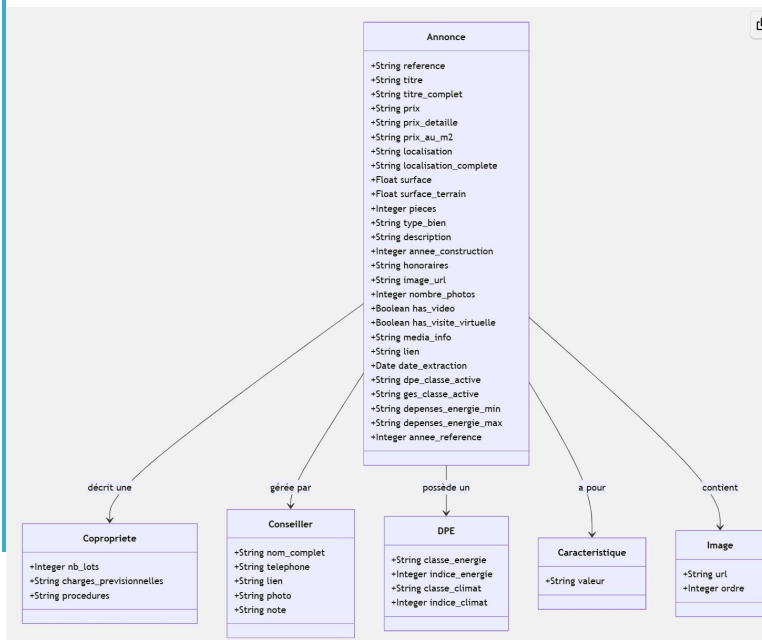
Un script Python robuste a été développé pour l'import et le nettoyage des données JSON, incluant la validation des prix, l'extraction des surfaces, la gestion des diagnostics DPE et le traitement des relations un-à-plusieurs. Le système gère automatiquement les doublons, les erreurs d'import et fournit des statistiques en temps réel, formant une solution ETL complète pour l'analyse du marché immobilier.

```
# Creation de la base de données imo_db
psql -U postgres -f sql\create_imo_db_init.sql

# Acces a imo_db et reation de la structure et des données partielles.
psql -U postgres -d imo_db -f sql\create_imo_db.sql

# Import des données dans la base imo_db a partir du fichier json
uv run import_annonces.py
```

Modèle UML de la base de données imo_sql



Compétence C5: Développement d'une API Rest

a APi Rest with Fast API from Postgres Datasource

services/c5-api

API Rest avec FASTAPI

Cette API Immobilière est développée avec FastAPI, un framework web moderne et haute performance pour Python 3.8+. L'architecture repose sur une base de données PostgreSQL avec une connexion gérée via psycopg2 et le curseur RealDictCursor pour un mapping naturel des résultats en dictionnaires Python.

L'application implémente le pattern Dependency Injection de FastAPI pour la gestion des connexions à la base de données, assurant une ouverture/fermeture propre des ressources à chaque requête.

Le système inclut un middleware CORS configuré pour autoriser les requêtes cross-origin, rendant l'API consommable par des applications frontend. La validation des paramètres de requête est assurée par le système de types de Pydantic intégré à FastAPI, avec des contraintes sur les valeurs numériques (bornes min/max) et la gestion des paramètres optionnels. L'API expose des endpoints RESTful conventionnels avec pagination, filtrage avancé (par type de bien, prix, surface, localisation) et recherche plein texte via des clauses ILIKE PostgreSQL.

La documentation interactive est automatiquement générée via Swagger UI et ReDoc, tandis que la configuration est externalisée via python-dotenv pour une gestion sécurisée des credentials. Les erreurs sont uniformément gérées avec des HTTPException standardisées et des codes d'état appropriés, garantissant une expérience client prédictible.

```
uv run api.py
```

Endpoints

GET	/	Root
GET	/annonces	Get Annonces
GET	/annonces/{annonce_id}	Get Annonce By Id
GET	/annonces/reference/{reference}	Get Annonce By Reference
GET	/annonces/search	Search Annonces
GET	/statistiques	Get Statistiques
GET	/types-bien	Get Types Bien
GET	/localisations	Get Localisations

Data Platform IMO-Ops pour un projet en intelligence artificielle dédié au marché Immobilier.

Dev IA RNCP 37827 / © 2025 Olivier LAVAUD

Conclusion

Ce projet de plateforme data immobilière IMO-OPS a permis la mise en œuvre d'une architecture complète et conteneurisée, orchestrant l'ensemble du cycle de vie des données.

L'infrastructure technique repose sur un écosystème Docker composé d'Airflow pour l'orchestration des pipelines, Redis est envisagée pour la gestion des files, PostgreSQL et MongoDB pour le stockage, ainsi que Grafana et des interfaces personnalisées pour la visualisation.

Cette stack technique robuste supporte une valeur métier significative, avec des calculs avancés de rentabilité, une segmentation fine du marché et des recommandations actionnables.

L'industrialisation sera garantie par l'automatisation des processus via les DAGs, la reproductibilité des environnements et le monitoring unifié, faisant de cette plateforme une solution opérationnelle et maintenable pour l'analyse du marché immobilier



imo_ops

Data Plateform IMO-Ops pour un projet en intelligence artificielle dédié au marché Immobilier.

Dev IA RNCP 37827 / © 2025 Olivier LAVAUD

Annexes

Documentation

pgAdmin 4

file Object Tools Edit View Window Help

Welcome imo_db/postgres@PostgreSQL 17* x postgres/postgres... x

imo_db/postgres@PostgreSQL 17

SQL

Showing rows: 1 to 90 Page No: 1 of 1

reference	titre	titre_coi	prix	prix_detaille	prix_au_m2	localisation	localisa	surface	surface_te	pieces	type_bien	description	
character varying (50)	character	character numeric (12,2)	character numeric (10,2)	character numeric (10,2)	character numeric (10,2)	character varying (255)	character numeric (8)	numeric (8)	numeric (8)	integer	character varying (10)	text	
1	1883607	Appa...	App...	189000.00	189 000 €	3436.36	Montpellier (34070)	App...	55.00	12.00	2	Appartement	En exclusivité – Montpellier, entre La Chamberte et La Martelle
2	1885312	Mais...	Mai...	315000.00	315 000 €	2581.97	MONTPELLIER (34080)	Mal...	122.00	250.00	4	Maison	*OPPORTUNITE *
3	1887434	Appa...	App...	169000.00	169 000 €	3840.91	Montpellier (34070)	App...	44.00	12.00	2	Appartement	APPARTEMENT T2 CLÉ EN MAIN - MONTPELLIER LES GRISETTES T
4	1883561	Appa...	App...	215000.00	215 000 €	4479.17	Montpellier (34090)	App...	48.00	[null]	3	Appartement	Appartement 3 Pièces Vue sur Parc - Calme Rare en Plein quartier de
5	1886031	Appa...	App...	145000.00	145 000 €	2132.35	Montpellier (34080)	App...	68.00	143.00	3	Appartement	APPARTEMENT T4 - 3 CHAMBRES - L'IMINIFIIX - TRAVERSANT - FON
6	1878714	Appa...	App...	379500.00	379 500 €	3872.45	Montpellier (34070)	App...	98.00	5.00	3	Appartement	
7	1886234	Appa...	App...	170000.00	170 000 €	3617.02	Montpellier (34070)	App...	47.00	12.00	2	Appartement	
8	1885348	Appa...	App...	145000.00	145 000 €	4027.78	Montpellier (34000)	App...	36.00	[null]	2	Appartement	
9	1874842	Appa...	App...	89000.00	89 000 €	3068.97	Montpellier (34080)	App...	29.00	[null]	2	Appartement	
10	1884128	Appa...	App...	169000.00	169 000 €	2600.00	Montpellier (34070)	App...	65.00	[null]	2	Appartement	
11	1882579	Studi...	Stud...	50000.00	50 000 €	2777.78	Montpellier (34090)	Stud...	18.00	[null]	2	Appartement	
12	1880084	Duple...	Dupl...	183000.00	183 000 €	2577.46	Montpellier (34080)	Dupl...	71.00	10.00	2	Appartement	
13	1881531	Appa...	App...	269000.00	269 000 €	3586.67	Montpellier (34070)	App...	75.00	[null]	3	Appartement	
14	1881990	Appa...	App...	186000.00	186 000 €	3152.54	Montpellier (34070)	App...	59.00	[null]	3	Appartement	
15	1880618	Appa...	App...	230000.00	230 000 €	3538.46	Montpellier (34070)	App...	65.00	5.00	3	Appartement	
16	1879922	Park...	Park...	13000.00	13 000 €	1181.82	Montpellier (34000)	Park...	11.00	[null]	2	Appartement	
17	1879821	Gara...	Gar...	15000.00	15 000 €	1500.00	Montpellier (34000)	Gar...	10.00	[null]	2	Appartement	
18	1878148	Park...	Park...	11900.00	11 900 €	1322.22	Montpellier (34070)	Park...	9.00	[null]	2	Appartement	
19	1877594	Gara...	Gar...	15000.00	15 000 €	1071.43	Montpellier (34080)	Gar...	14.00	[null]	2	Appartement	
20	1877544	Appa...	App...	205000.00	205 000 €	4183.67	Montpellier (34000)	App...	49.00	[null]	2	Appartement	

Compass

My Queries

CONNECTIONS (1)

Search connections

localhost:27017

admin

config

local

real_estate

properties

properties

localhost:27017 > real_estate > properties

Documents 19.6K Aggregations Schema Indexes 1 Validation

Type a query: { field: 'value' } or Generate query

EXPLAIN RESET FIND Options

25 1 - 25 of 19621

```
{
  "_id": ObjectId("672e199a5dbf4355e97d7780"),
  "id": 1,
  "address": "61 Boulevard du Montparnasse",
  "city": "Paris",
  "price": 498658.97,
  "size": 187.2,
  "rooms": 5,
  "type": "House",
  "transaction_date": "2019-12-08",
  "source": "ZSON",
  "extraction_date": "2024-11-08"
}
```

```
{
  "_id": ObjectId("672e199a5dbf4355e97d7781"),
  "id": 2,
  "address": "70 Rue du Faubourg Saint-Honoré",
  "city": "Lille",
  "price": 366651.87,
  "size": 78.36,
  "rooms": 2,
  "type": "House",
  "transaction_date": "2024-05-04",
  "source": "ZSON",
  "extraction_date": "2024-11-08"
}
```