

ERP/1: Аналітика

# Технічні вимоги

Локальна DWH-інфраструктура  
для аналізу та обробки великих даних

# Зміст

## Зміст

<b>1</b>	<b>Умовні скорочення та визначення</b>	<b>3</b>
<b>2</b>	<b>Загальні відомості</b>	<b>4</b>
2.1	Передумови . . . . .	4
2.2	Питання, що вирішуються . . . . .	4
<b>3</b>	<b>Призначення та цілі впровадження</b>	<b>4</b>
<b>4</b>	<b>Класифікація вимог</b>	<b>5</b>
4.1	Функціональні вимоги . . . . .	5
4.1.1	Вимоги до локального аналітичного сховища (DWH) та його архітектури .	5
4.1.2	Вимоги до інтеграції та ETL/ELT контурів . . . . .	5
4.1.3	Вимоги до векторизованого виконання запитів та обчислень . . . . .	5
4.2	Нефункціональні вимоги . . . . .	5
4.2.1	Апаратні вимоги до обчислювальних вузлів . . . . .	6
4.2.2	Вимоги до програмного стеку . . . . .	6
4.2.3	Вимоги до продуктивності та масштабування . . . . .	6
4.2.4	Вимоги до безпеки та захисту інформації (NIST SP 800-53) . . . . .	6
4.2.5	Вимоги до відмовостійкості та резервування . . . . .	6
<b>5</b>	<b>Додаток А. Порівняльний аналіз OLAP-платформ та засобів візуалізації</b>	<b>7</b>
5.1	Рекомендації щодо архітектурного вибору . . . . .	7

### Анотація

У цьому документі наведено детальні технічні вимоги до підсистеми ERP/1: Аналітика, яка розробляється як сучасний високопродуктивний сервіс локального сховища великих даних (DWH), аналітичних обчислень та побудови бізнес-звітів (BI). Платформа інтегрується з децентралізованою екосистемою ERP/1 та передбачає локальне виконання швидких OLAP-запитів без передачі конфіденційної інформації за межі установ. Архітектурні вимоги фокусуються на використанні C99-сумісних мінімалістичних технологій без застосування екосистеми Rust, з повною інтеграцією в Erlang/OTP середовище периферійних вузлів. Документ визначає апаратний стек, вимоги безпеки відповідно до КСЗІ та NIST SP 800-53, а також містить порівняльний аналіз цільових інструментів для вибору рішення.

# 1. Умовні скорочення та визначення

Терміни та скорочення, що використовуються в цьому документі:

Термін / Скорочення	Значення
DWH	Data Warehouse (Сховище даних)
OLAP	Online Analytical Processing (Аналітична обробка в реальному часі)
ETL	Extract, Transform, Load (Витягнення, перетворення, завантаження)
ELT	Extract, Load, Transform (Витягнення, завантаження, перетворення)
BI	Business Intelligence (Бізнес-аналітика та візуалізація)
NIF	Native Implemented Function (Вбудована Erlang-функція мовою C)
OTP	Open Telecom Platform (Програмна екосистема Erlang)
Parquet	Відкритий стовпчиковий формат зберігання даних
AVX-512	Розширення векторних інструкцій процесорів x86-64
KC3I	Комплексна система захисту інформації

## 2. Загальні відомості

### 2.1. Передумови

Аналіз та побудова звітів за великими масивами транзакційних даних, реєстрів та документів вимагає виділення окремого обчислювального контуру, оптимізованого під читання великої кількості записів. Водночас, обробка даних у державному та корпоративному секторах накладає такі суворі архітектурні обмеження:

1. **Конфіденційність та автономність:** Відповідно до вимог законодавства України щодо захисту персональних даних та державних реєстрів, аналітична обробка повинна відбуватися виключно в межах локального захищеного контуру (on-premises) установи без надсилання даних у хмари.
2. **Мінімізація залежностей стеку компіляції:** Для спрощення сертифікації за вимогами КСЗІ та розгортання на великій кількості периферійних пристроїв необхідно обмежити використання складних інструментаріїв збирання та компіляції (виключити Rust/Cargo), віддаючи перевагу перевіреним системним засобам C99.
3. **Інтеграція з Erlang/OTP:** - Аналітична підсистема повинна безшовно взаємодіяти з транзакційним контуром ERP/1, написаним на Erlang/Elixir, через нативні механізми портів або C-інтерфейси NIF.

### 2.2. Питання, що вирішуються

Впровадження модуля «Аналітика» вирішує такі ключові завдання:

- Усунення аналітичного навантаження (великих агрегаційних запитів) з транзакційної бази даних Mnesia/KVS.
- Швидкий аналіз історичних даних обсягом до десятків мільярдів записів безпосередньо на периферійних вузлах.
- Зберігання та обробка даних у стиснутому стовпчиковому форматі, що зменшує вимоги до дискового простору.
- Візуалізація звітів за допомогою інтегрованих інструментів бізнес-аналітики.

## 3. Призначення та цілі впровадження

Основною метою створення модуля є надання інструментів побудови аналітичних звітів, агрегації показників роботи установи (наприклад, статистика судових справ, фінансовий облік, обробка звернень) у реальному часі на локальних вузлах з використанням сумісних, легкозамінних та швидкодіючих відкритих технологій.

## 4. Класифікація вимог

### 4.1. Функціональні вимоги

#### 4.1.1. Вимоги до локального аналітичного сховища (DWH) та його архітектури

Система повинна реалізовувати локальне стовпчикове сховище даних на базі вбудованого аналітичного двигуна (DuckDB), що відповідає таким вимогам:

- **Стовпчикове зберігання та стиснення:** зберігання агрегованих таблиць та журналів подій у стовпчиковому форматі для максимізації швидкості читання із застосуванням алгоритмів стиснення (ZSTD, LZ4).
- **Сумісність форматів:** повна підтримка роботи з відкритими аналітичними форматами файлів (Parquet, CSV, JSON).
- **Самостійність та ізоляція:** функціонування у вигляді окремого самостійного процесу (або ізольованого системного демона) з власним пулом виділеної оперативної пам'яті та лімітами обчислювальних ядер CPU, що запобігає впливу аналітичних обчислень на основні транзакційні процеси.
- **Масштабованість:** підтримка горизонтального масштабування шляхом федеративного виконання SQL-запитів над множиною незалежних баз та розподілених Parquet-файлів, а також реалізація механізмів Peer-to-Peer реплікації даних між периферійними вузлами для забезпечення відмовостійкості.
- **Автономність:** сховище повинно повноцінно функціонувати в Air-Gapped режимі без підключення до зовнішніх мережевих сервісів.

#### 4.1.2. Вимоги до інтеграції та ETL/ELT контурів

Організація імпорту даних із транзакційної підсистеми ERP/1:

- Асинхронний експорт змін транзакційних таблиць Mnesia/KVS у аналітичне сховище.
- Батчеве завантаження (Batch Loading) історичних масивів даних у фоновому режимі з низьким пріоритетом використання процесора (nice).
- Можливість виконання легких ELT перетворень безпосередньо всередині сховища за допомогою декларативних SQL-запитів.

#### 4.1.3. Вимоги до векторизованого виконання запитів та обчислень

Аналітичний двигун повинен мати такі можливості:

- Векторизоване виконання запитів (обробка даних пакетами рядків, а не по одному), що дозволяє повністю задіяти кеш процесора.
- Використання багатопотокового паралелізму для обробки одного аналітичного SQL-запиту.
- Наявність C99-сумісного API інтерфейсу для низькорівневої інтеграції з кодом Erlang/Elixir.

## 4.2. Нефункціональні вимоги

### 4.2.1. Апаратні вимоги до обчислювальних вузлів

В якості апаратної платформи використовуються стандартні сервери або робочі станції канцелярій:

- Процесор: архітектура x86-64 з обов'язковою підтримкою векторних інструкцій AVX2 або AVX-512.
- Оперативна пам'ять: від 16 до 64 ГБ RAM (з оптимізацією під in-memory обчислення).
- Накопичувачі: NVMe SSD з високою швидкістю випадкового читання.

### 4.2.2. Вимоги до програмного стеку

Вимоги до технологічного стеку та системних обмежень:

- **Використання C99:** Усі низькорівневі бібліотеки, драйвери, Erlang NIF модулі та інтерфейсні обгортки повинні бути розроблені на чистому стандарті мови C99 без складних зовнішніх бібліотек C++.
- **Erlang/OTP екосистема:** Керування життєвим циклом DWH, оркестрація ETL-процесів та обробка аналітичних задач повинні виконуватися як стандартні OTP-додатки (через супервізори та gen\_servers).

### 4.2.3. Вимоги до продуктивності та масштабування

Система аналітики повинна забезпечувати такі метрики продуктивності:

- Швидкість сканування даних: не менше 100 мільйонів рядків на секунду на одному обчислювальному ядрі CPU.
- Час виконання складних агрегацій (GROUP BY, JOIN) на масивах до 100 млн записів: не більше 1.5 секунди.
- Можливість горизонтального масштабування шляхом об'єднання локальних периферійних баз у централізований аналітичний кластер.

### 4.2.4. Вимоги до безпеки та захисту інформації (NIST SP 800-53)

Захист аналітичних даних забезпечується такими засобами:

- **SC-28 (Protection of Information at Rest):** повне шифрування дисків за допомогою LUKS та інтеграції з TPM 2.0.
- **AC-3 (Access Enforcement):** ізоляція доступу на рівні рядків таблиць (Row-Level Security) та атрибутивний контроль прав (ABAC).
- **Network Isolation:** блокування зовнішнього трафіку за допомогою Network Policies у Cilium.

### 4.2.5. Вимоги до відмовостійкості та резервування

Забезпечення доступності сервісу аналітики:

- Підтримка механізму холодної та гарячої заміни вузлів без втрати історичних архівів.
- Асинхронне нічне резервне копіювання локальних даних у централізоване об'єктне сховище.

## 5. Додаток А. Порівняльний аналіз OLAP-платформ та засобів візуалізації

Для вибору цільового рішення нижче наведено порівняльний аналіз двох провідних відкритих аналітичних баз даних (ClickHouse та MonetDB як найлегшого конкурента) та системи візуалізації звітів (Apache Superset) під кутом архітектурних вимог ERP/1:

Критерій	ClickHouse	MonetDB	Apache Superset
Тип системи	Розподілена стовпчикова СКБД	Стовпчикова реляційна СКБД (standalone)	BI-платформа візуалізації та побудови дашбордів
Мова реалізації	C++20	Чистий C (C99)	Python (SQLAlchemy) + React (TypeScript)
Складність розгортання	Висока (вимагає серверного кластера, Zookeeper/Keeper)	Низька (легковагий демон без зовнішніх залежностей)	Середня (запускається як окремий контейнер у K8s)
Інтеграція Erlang/OTP	з Через TCP/HTTP клієнти або зовнішні драйвери	Через нативний Erlang-драйвер або C-порти	Через стандартні SQL-драйвери до ClickHouse/MonetDB
Цільове призначення	Централізований аналітичний хаб для мільярдів рядків	Автономне легке аналітичне сховище для локальних вузлів	Клієнтський інтерфейс бізнес-аналітики

### 5.1. Рекомендації щодо архітектурного вибору

1. Для децентралізованих периферійних вузлів (робочих станцій) рекомендується використання **DuckDB**. Проте, згідно з технічними вимогами, рішення на базі DuckDB має бути повністю **самостійним** (ізолювані системні ресурси, окремі пули потоків виконання NIF/Port) та **масштабованим** (з підтримкою розподіленого федеративного виконання запитів і Peer-to-Peer синхронізації Parquet-файлів), що дозволяє успішно конкурувати з важкими клієнт-серверними рішеннями на локальному рівні.
2. Для центрального аналітичного вузла системи ERP/1 рекомендується розгортання кластера **ClickHouse**. Він бере на себе роль глобального DWH, куди стікаються агреговані дані з усіх 600+ установок.
3. Для побудови та відображення звітів користувачам як окремий аналітичний додаток пропонується використовувати **Apache Superset**, який підключається як до локальних баз DuckDB, так і до центрального ClickHouse.