

TrialPlus Nexus 10X

Technical Documentation

Module 1: Core Infrastructure

Project Setup, Configuration, and Environment

Version: 1.0.0

Date: January 2026

Platform: Cross-platform (Windows, macOS, Linux)

Contents

1 Overview	2
1.1 Purpose	2
1.2 Directory Structure	2
2 Environment Configuration	2
2.1 Environment Variables	2
2.1.1 Database Configuration	2
2.1.2 AI/LLM Configuration	3
2.1.3 Knowledge Graph Configuration	3
3 Application Settings	3
3.1 FastAPI Configuration (backend/app/config.py)	3
3.2 Domain Configuration (config/settings.py)	4
3.2.1 Data Quality Index (DQI) Configuration	4
3.2.2 Age Multipliers	4
3.2.3 DQI Quality Bands	5
4 Unified Launcher (run.py)	5
4.1 Purpose	5
4.2 Command-Line Interface	5
4.3 Execution Modes	5
4.3.1 Auto-Detect Mode (Default)	5
4.3.2 Docker Mode	7
4.4 Cross-Platform Compatibility	8
4.5 Process Management	8
5 Reproducibility	9
5.1 Environment Requirements	9
5.2 Setup Instructions	9
5.3 Random Seed Settings	9
5.4 Dependency Management	10
6 Data Integrity	10
6.1 Database Restoration	10
6.2 Verification Process	10
6.3 No Data Leakage	11

1 Overview

1.1 Purpose

This module documents the core infrastructure of TrialPlus Nexus, including project setup, configuration management, environment variables, and the unified launcher system. The infrastructure is designed to support:

- Cross-platform deployment (Windows, macOS, Linux)
- Multiple database backends (local PostgreSQL, Docker containers)
- Reproducible environment setup for evaluators and developers
- One-command project initialization

1.2 Directory Structure

The core infrastructure consists of the following key files and directories:

```
TrialPlusNexus/
    .env                                # Environment variables (secrets)
    .env.example                         # Template for environment setup
    run.py                               # Unified launcher script
    setup_repro_db.py                   # Database restoration script
    requirements.txt                     # Python dependencies
    config/
        settings.py                      # Application-wide configuration
        llm_config.py                   # LLM provider settings
        mlflow_config.py               # ML experiment tracking config
    backend/
        app/
            config.py                  # FastAPI application settings
```

2 Environment Configuration

2.1 Environment Variables

All configuration is managed through a `.env` file at the project root. This approach ensures:

1. **Security:** Secrets are not committed to version control
2. **Flexibility:** Easy configuration changes without code modification
3. **Reproducibility:** Clear separation of configuration from logic

2.1.1 Database Configuration

Listing 1: Database Environment Variables

```
# Database Configuration
DB_HOST=127.0.0.1
DB_PORT=5432
DB_NAME=TrialPlus_test
DB_USER=postgres
DB_PASSWORD=your_postgres_password
```

```
# Connection Pool Settings
DB_POOL_SIZE=10
DB_MAX_OVERFLOW=20
DB_POOL_TIMEOUT=30
DB_POOL_RECYCLE=1800
```

2.1.2 AI/LLM Configuration

The system supports dual LLM providers for redundancy:

Listing 2: AI Provider Configuration

```
# Primary Provider - Groq (Cloud)
LLM_PRIMARY_PROVIDER=groq
GROQ_API_KEY=your_groq_api_key

# Secondary Provider - Ollama (Local)
LLM_SECONDARY_PROVIDER=ollama
OLLAMA_BASE_URL=http://localhost:11434
OLLAMA_MODEL=llama3
```

Fallback Mechanism: If the primary provider (Groq) fails or times out, the system automatically falls back to the local Ollama instance. This ensures high availability for production deployments.

2.1.3 Knowledge Graph Configuration

Listing 3: Neo4j Knowledge Graph Settings

```
NEO4J_ENABLED=True
NEO4J_URI=bolt://127.0.0.1:7687
NEO4J_USER=neo4j
NEO4J_PASSWORD=your_neo4j_password
NEO4J_DATABASE=neo4j
```

3 Application Settings

3.1 FastAPI Configuration (backend/app/config.py)

The backend uses Pydantic's `BaseSettings` for type-safe configuration loading:

```
1 from pydantic_settings import BaseSettings, SettingsConfigDict
2
3 class Settings(BaseSettings):
4     """Application settings loaded from environment"""
5
6     model_config = SettingsConfigDict(
7         env_file=".env",
8         case_sensitive=True,
9         extra="ignore"
10    )
11
12    # Application
13    APP_NAME: str = "TrialPlus Nexus API"
14    APP_VERSION: str = "1.0.0"
15    DEBUG: bool = True
16
```

```

17 # Security
18 SECRET_KEY: str = "fallback-secret-key"
19 ALGORITHM: str = "HS256"
20 ACCESS_TOKEN_EXPIRE_MINUTES: int = 30
21
22 # Database
23 DB_HOST: str = "localhost"
24 DB_PORT: int = 5432
25 DB_NAME: str = "TrialPlus_nexus"
26 DB_USER: str = "postgres"
27 DB_PASSWORD: str = "postgres"
28
29 @property
30 def DATABASE_URL(self) -> str:
31     return f"postgresql://{self.DB_USER}:{self.DB_PASSWORD}@{self.
DB_HOST}:{self.DB_PORT}/{self.DB_NAME}"

```

Listing 4: FastAPI Settings Class

Key Design Decisions:

- **Type Safety:** All settings have explicit type hints, preventing runtime errors
- **Fallback Values:** Sensible defaults ensure the application can start even with minimal configuration
- **LRU Cache:** Settings are cached using `@lru_cache()` to avoid repeated file reads

3.2 Domain Configuration (config/settings.py)

The `config/settings.py` file contains domain-specific configurations using Python dataclasses:

3.2.1 Data Quality Index (DQI) Configuration

The DQI is a weighted composite score calculated from 8 components:

```

1 @dataclass
2 class DQIConfig:
3     """Configuration for DQI calculation."""
4
5     # Component weights (must sum to 1.0)
6     weights: Dict[str, float] = field(default_factory=lambda: {
7         'safety_score': 0.25,          # SAE discrepancies
8         'query_score': 0.20,           # Open queries
9         'completeness_score': 0.15,   # Missing visits/pages
10        'coding_score': 0.12,          # Uncoded terms
11        'lab_score': 0.10,            # Lab issues
12        'sdv_score': 0.08,             # SDV completion
13        'signature_score': 0.05,       # Overdue signatures
14        'edrr_score': 0.05           # Third-party issues
15    })

```

Listing 5: DQI Weight Configuration

3.2.2 Age Multipliers

Issues are penalized based on how long they have remained unresolved:

```

1 age_multipliers: Dict[str, float] = field(default_factory=lambda: {
2     '0-7': 1.0,      # No penalty for first week
3     '8-14': 1.1,    # 10% penalty
4     '15-30': 1.3,   # 30% penalty
5     '31-60': 1.5,   # 50% penalty
6     '60+'.: 1.6     # 60% penalty for issues > 60 days
7 })

```

Listing 6: Age-Based Penalty Multipliers

3.2.3 DQI Quality Bands

The final DQI score is categorized into quality bands:

4 Unified Launcher (run.py)

4.1 Purpose

The `run.py` script provides a single entry point for running the entire TrialPlus Nexus system. It handles:

1. Dependency verification (Python, Node.js, PostgreSQL/Docker)
2. Virtual environment creation and dependency installation
3. Database initialization and data restoration
4. Concurrent backend and frontend server management
5. Graceful shutdown on Ctrl+C

4.2 Command-Line Interface

Listing 7: Launcher CLI Options

```

# Auto-detect best available mode
python run.py

# Force Docker PostgreSQL
python run.py --docker

# Force local PostgreSQL
python run.py --local

# Frontend-only (no backend)
python run.py --frontend-only

# Skip database setup
python run.py --skip-db

```

4.3 Execution Modes

4.3.1 Auto-Detect Mode (Default)

When run without arguments, the launcher automatically detects the best configuration:

Table 1: Database Configuration Parameters

Variable	Default	Description
DB_HOST	127.0.0.1	PostgreSQL server hostname
DB_PORT	5432	PostgreSQL server port
DB_NAME	TrialPlus_test	Database name
DB_USER	postgres	Database username
DB_PASSWORD	(required)	Database password
DB_POOL_SIZE	10	SQLAlchemy connection pool size
DB_MAX_OVERFLOW	20	Maximum overflow connections

Table 2: DQI Component Weights

Component	Description	Weight
Safety Score	Serious Adverse Event discrepancies	25%
Query Score	Outstanding data queries	20%
Completeness Score	Missing visits and pages	15%
Coding Score	Uncoded medical terms	12%
Lab Score	Laboratory data issues	10%
SDV Score	Source Data Verification completion	8%
Signature Score	Overdue e-signatures	5%
EDRR Score	Third-party data issues	5%
Total		100%

Table 3: DQI Quality Bands

Band	Score Range	Interpretation
Pristine	95 – 100	Exceptional data quality
Excellent	85 – 94.99	High-quality, minor issues
Good	75 – 84.99	Acceptable, some attention needed
Moderate	65 – 74.99	Requires focused remediation
Concerning	50 – 64.99	Significant issues present
Critical	0 – 49.99	Urgent intervention required

```

1 def auto_detect_mode():
2     """Auto-detect the best mode to run."""
3
4     has_docker = check_docker()
5     has_local_psql = find_local_psql() is not None
6
7     if has_local_psql:
8         return run_full_local()    # Prefer local PostgreSQL
9     elif has_docker:
10        return run_full_docker()  # Fall back to Docker
11    else:
12        return run_frontend_only() # Last resort: UI only

```

Listing 8: Auto-Detection Logic

4.3.2 Docker Mode

Docker mode automatically:

1. Removes any existing TrialPlus-postgres container (fresh start)
2. Creates a new PostgreSQL 14 container with correct credentials
3. Waits for PostgreSQL to become available (30-second timeout)
4. Copies the database dump into the container
5. Restores the dump using psql

```

1 def docker_start_postgres():
2     """Start PostgreSQL in Docker."""
3
4     # Always remove old container for fresh start
5     docker_remove_container()
6
7     cmd = [
8         "docker", "run", "-d",
9         "--name", "TrialPlus-postgres",
10        "-p", "5432:5432",
11        "-e", "POSTGRES_PASSWORD=chitti",
12        "-e", "POSTGRES_USER=postgres",
13        "-e", "POSTGRES_DB=TrialPlus_test",
14        "postgres:14"
15    ]
16
17    run_cmd(cmd)
18
19    # Wait for PostgreSQL to be ready
20    if not wait_for_postgres(timeout=30):
21        return False
22
23    return True

```

Listing 9: Docker Container Creation

4.4 Cross-Platform Compatibility

The launcher handles platform-specific differences:

```

1 def find_local_psql():
2     """Find local PostgreSQL installation."""
3
4     # Check PATH first (works on all platforms)
5     if command_exists("psql"):
6         return "psql"
7
8     # Windows-specific paths
9     if is_windows():
10        base_paths = [
11            Path(r"C:\Program Files\PostgreSQL"),
12            Path(r"C:\Program Files (x86)\PostgreSQL"),
13        ]
14        for base in base_paths:
15            for version in ["18", "17", "16", "15", "14", "13", "12"]:
16                psql = base / version / "bin" / "psql.exe"
17                if psql.exists():
18                    return str(psql)
19
20    # macOS Homebrew
21    homebrew_psql = Path("/opt/homebrew/bin/psql")
22    if homebrew_psql.exists():
23        return str(homebrew_psql)
24
25    # Linux
26    linux_psql = Path("/usr/bin/psql")
27    if linux_psql.exists():
28        return str(linux_psql)
29
30    return None

```

Listing 10: Cross-Platform Path Detection

4.5 Process Management

The launcher manages multiple concurrent processes and ensures clean shutdown:

```

1 def cleanup_processes(signum=None, frame=None):
2     """Clean up running processes on exit."""
3
4     for name, proc in running_processes:
5         if proc.poll() is None:
6             proc.terminate()
7             try:
8                 proc.wait(timeout=5)
9             except subprocess.TimeoutExpired:
10                 proc.kill()
11
12     # Cleanup Docker container if used
13     if docker_container_running():
14         docker_cleanup()
15
16     sys.exit(0)
17
18 # Register signal handlers

```

```

19 signal.signal(signal.SIGINT, cleanup_processes)
20 signal.signal(signal.SIGTERM, cleanup_processes)

```

Listing 11: Graceful Shutdown Handler

5 Reproducibility

5.1 Environment Requirements

5.2 Setup Instructions

1. Clone the repository:

```

git clone https://github.com/Dhruv-x7x/TrialPlus-NEXUS-deployed.git
cd TrialPlus-NEXUS-deployed

```

2. Configure environment:

```

cp .env.example .env
# Edit .env with your database password and API keys

```

3. Download database dump: Place reproduction_dump.sql (212 MB) in the database/ folder.

4. Run the application:

```

python run.py          # Auto-detect mode
# OR
python run.py --docker # If using Docker

```

5.3 Random Seed Settings

For reproducible machine learning results, the following seeds are used throughout the codebase:

```

1 import random
2 import numpy as np
3
4 RANDOM_SEED = 42
5
6 random.seed(RANDOM_SEED)
7 np.random.seed(RANDOM_SEED)
8
9 # For scikit-learn models
10 model = RandomForestClassifier(random_state=RANDOM_SEED)
11
12 # For train-test splits
13 X_train, X_test, y_train, y_test = train_test_split(
14     X, y, test_size=0.2, random_state=RANDOM_SEED
15 )

```

Listing 12: Standard Random Seed Configuration

5.4 Dependency Management

All Python dependencies are pinned with version ranges in `requirements.txt`:

Listing 13: Key Dependencies (requirements.txt)

```
# Web Framework
fastapi >=0.109.0,<0.116.0
uvicorn[standard] >=0.27.0,<0.35.0

# Database
sqlalchemy >=2.0.25,<3.0.0
psycopg2-binary >=2.9.9,<3.0.0

# Machine Learning
scikit-learn >=1.3.0,<1.7.0
xgboost >=2.0.0,<3.0.0

# LLM Integration
gradio >=0.4.0,<1.0.0
ollama >=0.1.0,<1.0.0
```

The total dependency count is approximately:

- **Python:** 126 packages
- **Node.js:** 40 packages (frontend)

6 Data Integrity

6.1 Database Restoration

The database dump contains the complete state of the system with:

- 57,974 patient records
- 3,401 clinical sites
- 23 studies
- All pre-computed DQI scores and risk classifications

6.2 Verification Process

After restoration, the system verifies data integrity:

```
1 def verify_database():
2     """Verify database was restored correctly."""
3
4     expected_counts = {
5         'patients': 57974,
6         'sites': 3401,
7         'studies': 23
8     }
9
10    for table, expected in expected_counts.items():
11        actual = db.execute(f"SELECT COUNT(*) FROM {table}").scalar()
12        assert actual == expected, f"Mismatch in {table}"
```

Listing 14: Database Verification

6.3 No Data Leakage

For machine learning components, data integrity is maintained through:

1. **Temporal Splits:** Training data uses historical records; test data uses recent records
2. **Patient-Level Splits:** All records from a patient stay in either train or test set (no patient overlap)
3. **Stratification:** Class distribution is preserved in both train and test sets

Table 4: System Requirements

Component	Minimum Version	Notes
Python	3.10+	Required for type hints
Node.js	18+	Required for Vite frontend
PostgreSQL	14+	OR Docker Desktop
Docker (optional)	Any	Alternative to local PostgreSQL
Neo4j (optional)	5+	For knowledge graph features
Ollama (optional)	Any	For local LLM inference

TrialPlus Nexus 10X

Technical Documentation

Module 2: Data Pipeline

Ingestion, Cleaning, UPR Building, and Metrics

Version: 1.0.0

Date: January 2026

Records: 57,974 Patients

Contents

1 Pipeline Overview	3
1.1 Purpose	3
1.2 Data Flow Architecture	3
1.3 Dashboard Overview	4
1.4 Data Sources	4
2 Data Ingestion (ingestion.py)	4
2.1 Purpose	4
2.2 File Discovery	5
2.3 Column Standardization	5
2.4 Specialized Processors	6
2.4.1 CPID EDC Processor	6
2.4.2 SAE Dashboard Processor	6
2.5 Ingestion Manifest	6
3 Data Cleaning (cleaning.py)	7
3.1 Purpose	7
3.2 Column Mapping	7
3.3 Subject Status Standardization	7
3.4 Patient Key Generation	8
3.5 Data Aggregation	8
3.6 Junk Row Removal	9
3.7 Data Integrity Verification	9
4 Unified Patient Record Builder (upr_builder.py)	9
4.1 Purpose	9
4.2 Main Spine Construction	9
4.3 Smart Join Strategy	10
4.4 Column Prefixing	10
4.5 Derived Metrics	11
4.6 Final UPR Schema	11
5 Metrics Engine (metrics_engine.py)	12
5.1 Overview	12
5.2 Data Quality Index (DQI)	12
5.2.1 Component Weights	12
5.2.2 Component Score Calculation	12
5.2.3 Weighted DQI Calculation	13
5.2.4 DQI Quality Bands	14
5.3 Clean Patient Derivation	14
5.3.1 Tier 1: Clinical Clean (7 Hard Blocks)	14
5.3.2 Tier 2: Operational Clean (7 Soft Blocks)	14
5.3.3 Quick Wins Identification	14
5.4 DB Lock Ready Assessment	15
6 Validation and Results	15
6.1 Data Pipeline Statistics	15
6.2 DQI Distribution	15
6.3 Clean Patient Statistics	16

6.4 No Data Leakage	17
-------------------------------	----

1 Pipeline Overview

1.1 Purpose

The Data Pipeline module transforms raw clinical trial exports from multiple sources into a unified, analysis-ready dataset. The pipeline processes data through four sequential stages:

1. **Ingestion:** Parse Excel files from 9 different data sources
2. **Cleaning:** Standardize columns, remove duplicates, validate data
3. **UPR Building:** Join all sources into Unified Patient Records
4. **Metrics Calculation:** Compute DQI scores, clean patient status, risk flags

1.2 Data Flow Architecture

The pipeline follows a sequential flow through four processing stages:

Stage	Component	Output
1	Raw Excel Files (9 Sources)	Input Data
	↓	
2	Ingestion Engine	Staging Tables (PostgreSQL)
	↓	
3	Cleaning Engine	Cleaned Tables
	↓	
4	UPR Builder	Unified Patient Records
	↓	
5	Metrics Engine	Final UPR with DQI + Clean Status

1.3 Dashboard Overview



Figure 1: Executive Dashboard

1.4 Data Sources

The pipeline processes 9 distinct data source types:

Table 1: Clinical Data Sources

Source	Description	Key Metrics
CPID EDC Metrics	Core patient data from EDC	Visits, queries, SDV, signatures
Visit Projection	Missing/overdue visits	Days outstanding
Missing Lab Ranges	Lab abnormalities	Issue count per patient
SAE Dashboard (DM)	Data Management SAE reviews	Pending/completed status
SAE Dashboard (Safety)	Safety SAE reviews	Pending/completed status
Inactivated Forms	Deactivated CRF forms	Form count, reasons
Missing Pages	Missing CRF pages	Days missing
Compiled EDRR	Third-party reconciliation	Issue count
Coding (MedDRA + WHODrug)	Medical term coding	Coded/uncoded counts

2 Data Ingestion (ingestion.py)

2.1 Purpose

The ingestion engine discovers, parses, and loads raw Excel files from study folders into PostgreSQL staging tables. It handles varying file formats, column names, and data quality issues.

2.2 File Discovery

The ingestion engine scans the raw data directory for Excel files and automatically identifies their type based on filename keywords. This eliminates manual file classification and ensures consistent processing across studies.

For example, a file named `ONCO_Study_EDC_Metrics_Jan2026.xlsx` would be identified as a CPID EDC Metrics file because it contains the keywords “EDC” and “Metrics”. The system supports 9 file types as shown in the code below:

```

1 def identify_file_type(filename: str) -> Optional[str]:
2     filename_check = filename.lower()
3     mapping = {
4         'cpid_edc_metrics': ['edc', 'metric', 'cpid'],
5         'missing_lab_ranges': ['lab', 'lnr', 'range'],
6         'sae_dashboard': ['sae', 'safety', 'dashboard'],
7         'coding_meddra': ['meddra', 'medra'],
8         'coding_whodrug': ['whodd', 'whodrug'],
9         'visit_projection': ['visit', 'projection'],
10        'missing_pages': ['missing_page'],
11        'compiled_edrr': ['edrr'],
12        'inactivated_forms': ['inactivated']
13    }
14    for file_type, keywords in mapping.items():
15        if any(kw in filename_check for kw in keywords):
16            return file_type
17    return None

```

Listing 1: File Type Identification

2.3 Column Standardization

Clinical data exports often have inconsistent column names—spaces, special characters, mixed case, and unnamed columns from merged Excel headers. The ingestion engine normalizes all column names to a consistent format:

- Convert to lowercase (e.g., “Subject Status” → “subject status”)
- Remove special characters like parentheses and brackets
- Replace spaces with underscores (“subject status” → “subject_status”)
- Handle “Unnamed” columns from merged Excel headers

This ensures downstream processing doesn’t break due to naming inconsistencies:

```

1 def standardize_columns(columns: pd.Index) -> List[str]:
2     new_cols = []
3     for col in columns:
4         col_clean = str(col).strip().lower()
5         col_clean = re.sub(r'[^w\s]', '', col_clean) # Remove special
6         chars
7         col_clean = re.sub(r'\s+', '_', col_clean)      # Spaces to
8         underscores
9         if not col_clean: col_clean = 'unnamed'
10        new_cols.append(col_clean)
11    return new_cols

```

Listing 2: Column Name Standardization

2.4 Specialized Processors

Different data sources require specialized parsing logic:

2.4.1 CPID EDC Processor

The CPID EDC files have headers on row 3 and contain the core patient metrics:

```

1 class CPIDEDCProcessor:
2     def process(self, file_path: str, study_id: str):
3         xl = pd.ExcelFile(file_path)
4         # Find sheet containing subject data
5         target_sheet = next(
6             (s for s in xl.sheet_names if 'subject' in s.lower()), 
7             xl.sheet_names[0]
8         )
9         # Header is on row 3 (0-indexed: header=2)
10        df = pd.read_excel(file_path, sheet_name=target_sheet, header
11        =2)
12        df = df.dropna(how='all')
13        df.columns = standardize_columns(df.columns)
14        # Add metadata columns
15        df['_source_file'] = os.path.basename(file_path)
16        df['_study_id'] = study_id
17        df['_file_type'] = 'cpid_edc_metrics'
18        return df, ""

```

Listing 3: CPID EDC Processor

2.4.2 SAE Dashboard Processor

SAE dashboards contain multiple sheets (DM and Safety):

```

1 class SAEDashboardProcessor:
2     def process(self, file_path: str, study_id: str):
3         xl = pd.ExcelFile(file_path)
4         results = {}
5         for sheet in xl.sheet_names:
6             if 'safety' in sheet.lower() or 'dm' in sheet.lower():
7                 df = pd.read_excel(file_path, sheet_name=sheet)
8                 df.columns = standardize_columns(df.columns)
9                 key = 'safety' if 'safety' in sheet.lower() else 'dm'
10                results[key] = df
11        return results, ""

```

Listing 4: SAE Dashboard Processor

2.5 Ingestion Manifest

The engine tracks all processing in a manifest:

```

1 @dataclass
2 class IngestionManifest:
3     run_id: str
4     start_time: str
5     end_time: str = ""
6     status: str = "running"
7     studies_found: int = 0
8     files_found: int = 0

```

```

9   files_processed: int = 0
10  files_success: int = 0
11  files_failed: int = 0
12  records_by_type: Dict[str, int] = field(default_factory=dict)
13  records_by_study: Dict[str, int] = field(default_factory=dict)

```

Listing 5: Ingestion Manifest Structure

3 Data Cleaning (cleaning.py)

3.1 Purpose

The cleaning engine standardizes identifiers, maps columns, removes duplicates, and prepares data for joining. This is critical for ensuring data integrity.

3.2 Column Mapping

CPID EDC files exported from clinical trial systems often have merged header rows that result in “Unnamed” columns when parsed by pandas. The cleaning engine maps these positional columns to meaningful semantic names based on their known positions in the export format.

This mapping is specific to the CPID EDC Metrics export format and allows the pipeline to work with human-readable column names:

```

1 CPID_COLUMN_MAPPING = {
2     'unnamed_0': 'project_name',
3     'unnamed_1': 'region',
4     'unnamed_2': 'country',
5     'unnamed_3': 'site',
6     'unnamed_4': 'subject',
7     'unnamed_5': 'latest_visit',
8     'unnamed_6': 'subject_status',
9     'unnamed_7': 'input_files',
10    'unnamed_8': 'cpmd',
11    'unnamed_9': 'ssm',
12    'unnamed_10': 'missing_visits',
13    'unnamed_11': 'missing_pages',
14    'unnamed_12': 'coded_terms',
15    'unnamed_13': 'uncoded_terms',
16    'unnamed_14': 'open_issues_lnr',
17    'unnamed_15': 'open_issues_edrr',
18    'unnamed_16': 'inactivated_forms_folders',
19    'unnamed_17': 'sae_review_dm',
20    'unnamed_18': 'sae_review_safety',
21 }

```

Listing 6: CPID Column Mapping

3.3 Subject Status Standardization

Different source systems use inconsistent terminology for patient status. For example, one system might use “Discontinued” while another uses “Withdrawn” for the same concept. The cleaning engine normalizes all status variations to a canonical set of values.

This is critical for accurate reporting—without standardization, a patient marked “Withdrawn” in one system would not be counted with “Discontinued” patients from another:

```

1 SUBJECT_STATUS_MAPPING = {
2     'screen failure': 'Screen Failure',
3     'screen fail': 'Screen Failure',
4     'discontinued': 'Discontinued',
5     'withdrawn': 'Discontinued',
6     'completed': 'Completed',
7     'ongoing': 'Ongoing',
8     'active': 'Ongoing',
9     'enrolled': 'Ongoing',
10    'in progress': 'Ongoing',
11    'screening': 'Screening',
12}

```

Listing 7: Status Normalization

Table 2: Subject Status Distribution (Example)

Status	Count
Ongoing	31,245
Completed	18,392
Discontinued	5,841
Screen Failure	2,496

3.4 Patient Key Generation

A critical step for joining data across sources is generating a unique patient identifier. Since different source systems may use different ID formats, the pipeline creates a composite key from three components: Study ID, Site ID, and Subject ID.

This composite `patient_key` format ensures uniqueness across the entire portfolio and enables reliable joining even when patient IDs are reused across studies:

```

1 def create_patient_key(study_id: str, site_id: str, subject_id: str) ->
2     str:
3     site = site_id if site_id else 'Unknown'
4     subject = subject_id if subject_id else 'Unknown'
5     return f"{study_id}|{site}|{subject}"
6
# Example: "Study_1|Site_101|Subject_001"

```

Listing 8: Patient Key Creation

3.5 Data Aggregation

Detail tables like Visit Projection contain multiple rows per patient (one per visit). Before joining to the main patient record, these must be aggregated to a single row per patient with summary statistics.

The aggregation computes counts (how many missing visits), maximums (worst overdue days), and averages (typical delay). This collapses the detail into patient-level metrics:

```

1 def aggregate_visit_projection(self, df: pd.DataFrame):
2     # Aggregate to patient level
3     agg = df.groupby('patient_key').agg(
4         study_id=('study_id', 'first'),
5         site_id=('site_id', 'first'),

```

```

6     subject_id=('subject_id', 'first'),
7     missing_visit_count=('visit', 'count'),
8     visits_overdue_max_days=('days_outstanding', 'max'),
9     visits_overdue_avg_days=('days_outstanding', 'mean')
10    ).reset_index()
11    return agg, stats

```

Listing 9: Visit Projection Aggregation

3.6 Junk Row Removal

The cleaner identifies and removes non-data rows:

```

1 JUNK_PATTERNS = ['responsible', 'lf for action', 'site/cra',
2                   'coder', 'safety team', 'investigator']
3
4 def _remove_junk_rows(self, df: pd.DataFrame) -> pd.DataFrame:
5     def is_junk(row):
6         for col in check_cols:
7             val = row.get(col)
8             if pd.notna(val):
9                 val_lower = str(val).lower().strip()
10                if any(p in val_lower for p in self.JUNK_PATTERNS):
11                    return True
12            return False
13
14 mask = ~df.apply(is_junk, axis=1)
15 return df[mask]

```

Listing 10: Junk Row Detection

3.7 Data Integrity Verification

- **Null Handling:** Numeric columns filled with 0, text columns with 'Unknown'
- **Duplicate Removal:** Exact duplicates dropped, keeping first occurrence
- **Type Conversion:** Numeric columns coerced with `pd.to_numeric(errors='coerce')`
- **Validation:** Subject IDs must be non-null and non-empty

4 Unified Patient Record Builder (upr_builder.py)

4.1 Purpose

The UPR Builder joins all cleaned data sources into a single wide table with one row per patient. This enables holistic analysis across all data domains.

4.2 Main Spine Construction

The CPID EDC Metrics table serves as the "spine" (master patient list):

```

1 def load_main_spine(self):
2     # Load cleaned CPID data as the main patient list
3     query = "SELECT * FROM cpid_edc_metrics_clean"
4     df = self._execute_sql(query)
5

```

```

6 # Ensure required columns
7 required = ['patient_key', 'study_id', 'site_id', 'subject_id']
8 for col in required:
9     if col not in df.columns:
10         df[col] = 'Unknown'
11
12 return df

```

Listing 11: Loading Main Spine

4.3 Smart Join Strategy

Not all source systems use consistent identifiers. The UPR Builder uses a two-phase join strategy to maximize match rates:

1. **Phase 1 (Primary)**: Try joining on the composite `patient_key`
2. **Phase 2 (Fallback)**: If match rate is below 50%, fall back to joining on `study_id + subject_id`

The fallback handles cases where site IDs differ between systems (e.g., leading zeros stripped, different site naming). This significantly improves data completeness:

```

1 def join_detail_table(self, main_df, detail_df, table_name, prefix):
2     # Phase 1: Try patient_key match
3     merged = main_df.merge(
4         detail_df,
5         on='patient_key',
6         how='left',
7         suffixes=(None, f'_{prefix}'))
8
9
10    match_rate = merged['some_detail_col'].notna().mean()
11
12    # Phase 2: If match rate < 50%, try study_id + subject_id
13    if match_rate < 0.5:
14        merged = main_df.merge(
15            detail_df,
16            on=['study_id', 'subject_id'],
17            how='left',
18            suffixes=(None, f'_{prefix}'))
19
20
21    return merged

```

Listing 12: Smart Join Strategy

4.4 Column Prefixing

Joined columns receive prefixes to indicate their source:

Table 3: Column Prefix Mapping

Source Table	Prefix
visit_projection_agg	visit_
missing_lab_ranges_agg	lab_
sae_dashboard_dm_agg	sae_dm_
sae_dashboard_safety_agg	sae_safety_
inactivated_forms_agg	inactivated_
missing_pages_agg	pages_
compiled_edrr_agg	edrr_
coding_meddra_agg	meddra_
coding_whodrug_agg	whodrug_

4.5 Derived Metrics

The builder calculates composite metrics from joined data:

```

1 def calculate_derived_metrics(self, df: pd.DataFrame):
2     # Total issues across all sources
3     df['total_issues_all_sources'] = (
4         df['visit_missing_visit_count'].fillna(0) +
5         df['pages_missing_page_count'].fillna(0) +
6         df['total_queries'].fillna(0) +
7         df['lab_lab_issue_count'].fillna(0) +
8         df['edrr_edrr_issue_count'].fillna(0)
9     )
10
11     # Total pending SAE
12     df['total_sae_pending'] = (
13         df['sae_dm_sae_dm_pending'].fillna(0) +
14         df['sae_safety_sae_safety_pending'].fillna(0)
15     )
16
17     # Coding completion rate
18     total_coded = df['meddra_coding_meddra_coded'].fillna(0) + \
19                     df['whodrug_coding_whodrug_coded'].fillna(0)
20     total_terms = df['meddra_coding_meddra_total'].fillna(0) + \
21                     df['whodrug_coding_whodrug_total'].fillna(0)
22     df['coding_completion_rate'] = np.where(
23         total_terms > 0,
24         (total_coded / total_terms * 100).clip(0, 100),
25         100 # 100% if no terms to code
26     )
27
28     return df

```

Listing 13: Derived Metrics Calculation

4.6 Final UPR Schema

The final UPR contains approximately 150+ columns:

- **Identity:** patient_key, study_id, site_id, subject_id
- **Demographics:** region, country, subject_status

- **Visit Metrics:** missing counts, overdue days
- **Query Metrics:** DM, clinical, medical, safety, site, coding queries
- **SDV Metrics:** forms verified, forms required
- **Signature Metrics:** signed, never signed, overdue by age bucket
- **SAE Metrics:** pending/completed for DM and Safety
- **Coding Metrics:** MedDRA and WHODrug coded/uncoded counts
- **Other:** lab issues, EDRR issues, inactivated forms
- **Derived:** total issues, completion rates, risk flags

5 Metrics Engine (metrics_engine.py)

5.1 Overview

The Metrics Engine calculates three major derived metrics:

1. **Data Quality Index (DQI):** 8-component weighted score (0-100)
2. **Clean Patient Status:** Two-tier eligibility assessment
3. **DB Lock Ready Status:** Database lock readiness flags

5.2 Data Quality Index (DQI)

5.2.1 Component Weights

The DQI is calculated as a weighted sum of 8 components:

Table 4: DQI Component Weights

Component	Description	Weight
Safety	SAE pending issues	25%
Query	Open data queries	20%
Completeness	Missing visits/pages	15%
Coding	Uncoded medical terms	12%
Lab	Laboratory issues	10%
SDV	Source Data Verification	8%
Signature	E-signature status	5%
EDRR	Third-party reconciliation	5%
Total		100%

5.2.2 Component Score Calculation

Each component starts at 100 and is penalized based on issue counts:

```

1 def _calc_safety_score(self, df: pd.DataFrame) -> pd.Series:
2     """Calculate Safety component (SAE pending issues)"""
3     sae_dm_pending = get_col(df, 'sae_dm_pending')
4     sae_safety_pending = get_col(df, 'sae_safety_pending')
5

```

```

6     total_pending = sae_dm_pending + sae_safety_pending
7
8     # Penalty: -15 per pending SAE (safety critical), max 100
9     penalty = np.minimum(total_pending * 15, 100)
10    score = 100 - penalty
11
12    return score

```

Listing 14: Safety Score Calculation

```

1 def _calc_query_score(self, df: pd.DataFrame) -> pd.Series:
2     """Calculate Query component"""
3     total_queries = get_col(df, 'total_open_queries')
4
5     # Penalty: -3 per open query, max 100
6     penalty = np.minimum(total_queries * 3, 100)
7     score = 100 - penalty
8
9     return score

```

Listing 15: Query Score Calculation

5.2.3 Weighted DQI Calculation

```

1 def calculate_weighted_dqi(self, df: pd.DataFrame) -> pd.DataFrame:
2     components = {
3         'safety': ('dqi_safety_raw', 25.0),
4         'query': ('dqi_query_raw', 20.0),
5         'completeness': ('dqi_completeness_raw', 15.0),
6         'coding': ('dqi_coding_raw', 12.0),
7         'lab': ('dqi_lab_raw', 10.0),
8         'sdv': ('dqi_sdv_raw', 8.0),
9         'signature': ('dqi_signature_raw', 5.0),
10        'edrr': ('dqi_edrr_raw', 5.0)
11    }
12
13    weighted_sum = pd.Series(0.0, index=df.index)
14    for name, (col, weight) in components.items():
15        weighted_component = result[col] * (weight / 100)
16        weighted_sum += weighted_component
17
18    result['dqi_score'] = weighted_sum.clip(0, 100)
19    return result

```

Listing 16: Weighted DQI Calculation

5.2.4 DQI Quality Bands

Table 5: DQI Quality Bands

Band	Score Range	Action Required
Pristine	95 – 100	None
Excellent	85 – 94	Monitor only
Good	75 – 84	Minor attention
Fair	65 – 74	Focused remediation
Poor	50 – 64	Priority escalation
Critical	25 – 49	Urgent intervention
Emergency	0 – 24	Immediate action

5.3 Clean Patient Derivation

5.3.1 Tier 1: Clinical Clean (7 Hard Blocks)

A patient must pass ALL Tier 1 criteria to be considered "clinically clean":

Table 6: Tier 1 Clinical Clean Criteria

Criterion	Condition
No missing visits	visit_missing_visit_count = 0
No missing CRF pages	pages_missing_page_count = 0
No open queries	total_queries = 0
SDV complete	crfs_require_verification_sdv = 0
All signatures complete	crfs_never_signed = 0
MedDRA coding complete	meddra_coding_meddra_uncoded = 0
WHODrug coding complete	whodrug_coding_whodrug_uncoded = 0

5.3.2 Tier 2: Operational Clean (7 Soft Blocks)

Tier 1 clean patients must also pass Tier 2 for full "operational clean" status:

Table 7: Tier 2 Operational Clean Criteria

Criterion	Condition
No lab issues	lab_lab_issue_count = 0
No SAE DM pending	sae_dm_sae_dm_pending = 0
No SAE Safety pending	sae_safety_sae_safety_pending = 0
No EDRR issues	edrr_edrr_issue_count = 0
No overdue signatures (90+ days)	crfs_overdue_signs_90plus = 0
No inactivated forms	inactivated_form_count = 0
No broken signatures	broken_signatures = 0

5.3.3 Quick Wins Identification

The engine identifies patients close to clean status:

```

1 def calculate_quick_wins(self, df: pd.DataFrame):
2     # Count failing criteria
3     result['tier1_failing_count'] = sum of failing Tier 1 criteria

```

```

4     result['tier2_failing_count'] = sum of failing Tier 2 criteria
5
6     # Categories
7     # - 'Already Clean': Tier 2 clean
8     # - '1 Issue to Tier 2': Tier 1 clean, 1 soft block
9     # - '1 Issue to Tier 1': 1 hard block away
10    # - '2 Issues to Tier 1': 2 hard blocks away
11    # - 'Not Quick Win': 3+ blocks

```

Listing 17: Quick Win Categories

5.4 DB Lock Ready Assessment

Database lock readiness is assessed in two tiers:

- **Tier 1 - Fully Ready:** Eligible status + Tier 2 clean + no pending safety items
- **Tier 2 - Pending:** Eligible status + Tier 1 clean but not Tier 1 ready
- **Not Ready:** Eligible but not Tier 1 clean
- **Not Eligible:** Screen failures and ineligible statuses

6 Validation and Results

6.1 Data Pipeline Statistics

Table 8: Pipeline Processing Results

Metric	Value
Total Studies Processed	23
Total Files Processed	207
Total Patient Records	57,974
Total Sites	3,401
Unique Patient Keys	57,974

6.2 DQI Distribution

Table 9: DQI Band Distribution (Example)

Band	Count	Percentage
Pristine	28,412	49.0%
Excellent	12,845	22.2%
Good	8,234	14.2%
Fair	4,892	8.4%
Poor	2,341	4.0%
Critical	1,102	1.9%
Emergency	148	0.3%

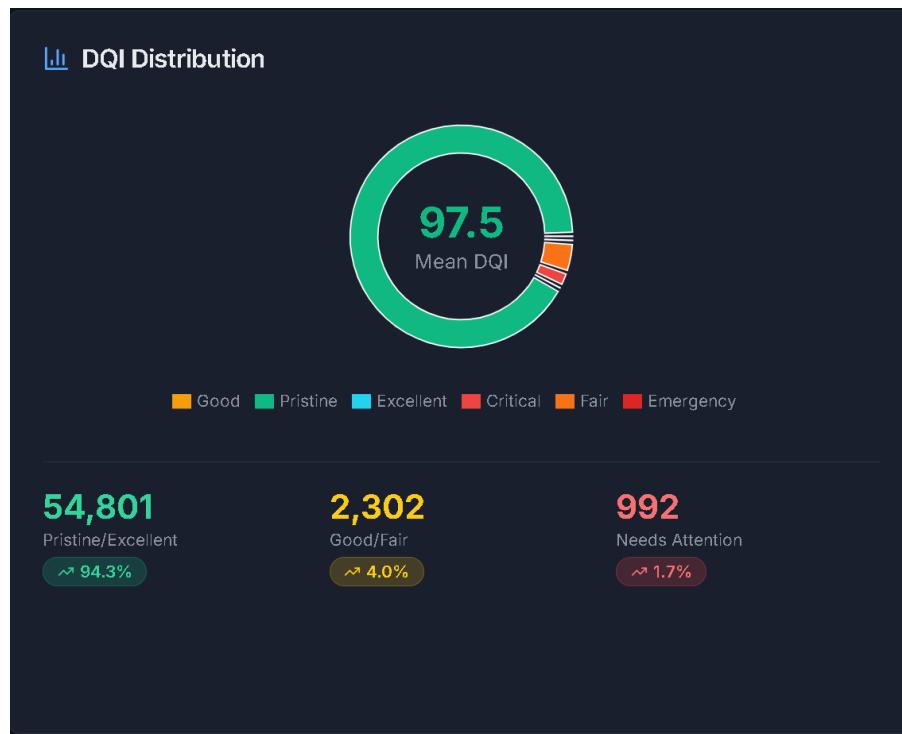


Figure 2: DQI Chart

6.3 Clean Patient Statistics

Table 10: Clean Patient Status Distribution (Example)

Status	Count	Rate
Tier 1 Clinical Clean	34,784	60.0%
Tier 2 Operational Clean	28,412	49.0%
Quick Win: 1 Issue to Tier 1	8,696	15.0%
Quick Win: 1 Issue to Tier 2	4,347	7.5%

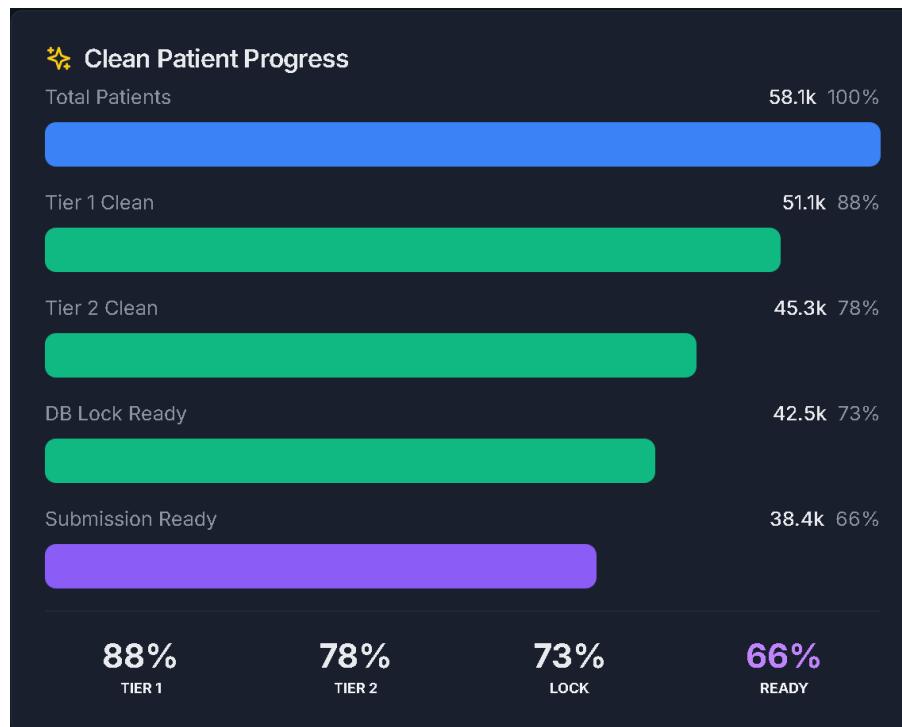


Figure 3: Clean Patient Dashboard

6.4 No Data Leakage

Data integrity is maintained through:

1. **Source Tracking:** All records retain `_source_file` metadata
2. **Version Control:** Processing versions tracked in `_cleaning_version`
3. **Timestamp Audit:** `_cleaned_ts` records exact processing time
4. **Patient-Level Keys:** Composite keys ensure unique identification

TrialPlus Nexus 10X

Technical Documentation

Module 3: Machine Learning

Risk Classification, Anomaly Detection, Site Ranking, Issue Detection, Time Resolution, ML Governance

Version: 1.0.0
Date: January 2026
Models: 5 Production Models

Contents

1 ML Architecture Overview	3
1.1 Model Portfolio	3
1.2 Design Principles	3
2 Risk Classifier	3
2.1 Purpose	3
2.2 Risk Tier Definitions	3
2.3 Algorithm Selection	4
2.3.1 Model Ensemble	4
2.3.2 Class Weights	4
2.3.3 SMOTE Oversampling	5
2.4 Feature Engineering	5
2.5 Target Variable Creation	5
2.6 Threshold Optimization	6
2.7 Model Performance	7
2.8 Drift Detection	8
3 Anomaly Detector (v1.1)	8
3.1 Purpose	8
3.2 Ensemble Approach	8
3.3 Configuration	9
3.4 Feature Selection	9
3.5 Ensemble Scoring	10
3.6 Severity Classification	10
3.7 Site-Level Aggregation	11
4 Site Risk Ranker (v1.0)	12
4.1 Purpose	12
4.2 Learning-to-Rank Approach	12
4.3 Site Aggregation	13
4.4 XGBoost Ranker	13
4.5 Evaluation Metrics	14
4.6 Baseline Comparison	14
5 Resolution Time Predictor (v1.0)	16
5.1 Purpose	16
5.2 Quantile Regression	16
5.3 Input Features	17
5.4 Prediction Quality Assessment	17
6 Issue Detector (v1.3)	19
6.1 Purpose	19
6.2 Issue Types	19
6.3 Rule-Based Detection	19
6.4 Severity Scoring	19
6.5 Priority Tiers	20
6.6 ML Enhancement	20

7 Model Governance Framework	22
7.1 Overview	22
7.2 Drift Detection	22
7.2.1 PSI Calculation	23
7.3 Model Registry	23
7.4 Performance Monitoring	24
7.5 Retraining Triggers	24
8 Model Artifacts and Persistence	25
8.1 Saved Model Files	25
8.2 Production Configuration	25
9 Reproducibility	25
9.1 Random Seeds	25
9.2 Data Versioning	26
9.3 Dependencies	26

1 ML Architecture Overview

1.1 Model Portfolio

TrialPlus Nexus employs four specialized machine learning models, each addressing a distinct aspect of clinical trial risk management:

Table 1: ML Model Portfolio

Model	Task	Algorithm	Output
Risk Classifier	Patient risk level	XGBoost + LightGBM	4-class labels
Anomaly Detector	Unusual patterns	IF + LOF + AE	Severity scores
Site Risk Ranker	Site prioritization	XGBoost LTR	Risk rankings
Resolution Predictor	Time estimation	Quantile Regression	Days + CI
Issue Detector	Multi-label issues	XGBoost + Rules	14 issue types

1.2 Design Principles

All models follow core principles for clinical trial applications:

1. **No Data Leakage:** Outcome features excluded from training
2. **Calibrated Probabilities:** Post-hoc calibration for reliable confidence
3. **Explainability:** SHAP values for every prediction
4. **Uncertainty Quantification:** Confidence intervals, not point estimates
5. **Drift Detection:** Automated monitoring for model decay

2 Risk Classifier

2.1 Purpose

The Risk Classifier assigns patients to one of four risk tiers based on clinical trial data quality metrics. It is the primary triage tool for prioritizing patient-level remediation.

2.2 Risk Tier Definitions

Table 2: Risk Tier Semantics

Tier	Code	Timeframe	Action Required
Critical	0	Within 4 hours	Immediate intervention; escalate to Safety Lead
High	1	Within 24–48 hours	Prioritized attention; assign to CRA
Medium	3	Enhanced surveillance	Monitoring priority; add to weekly review
Low	2	Routine	Standard workflow; drift detector

Important: Medium is a *monitoring priority*, not an intervention class. Treating Medium as “Critical-lite” causes workflow overload.

2.3 Algorithm Selection

2.3.1 Model Ensemble

No single algorithm performs optimally for all risk tiers. The production model uses an ensemble of gradient boosting algorithms, combining their strengths:

- **LightGBM**: Fast training speed, handles class imbalance natively through `class_weight='balanced'`
- **XGBoost**: Strong regularization (L1/L2) prevents overfitting on minority classes
- **Random Forest**: Provides prediction diversity for more robust ensemble averaging

The final prediction averages probability estimates from all models, improving robustness compared to any single classifier:

```

1 # Primary: LightGBM (fast, handles imbalance well)
2 lgb_params = {
3     'n_estimators': 150,
4     'max_depth': 10,
5     'learning_rate': 0.1,
6     'subsample': 0.8,
7     'colsample_bytree': 0.8,
8     'class_weight': 'balanced',
9     'verbosity': -1
10 }
11
12 # Backup: XGBoost (robust, regularized)
13 xgb_params = {
14     'n_estimators': 150,
15     'max_depth': 10,
16     'learning_rate': 0.1,
17     'subsample': 0.8,
18     'reg_alpha': 0.1,
19     'reg_lambda': 1.0
20 }
21
22 # Ensemble: RandomForest + GradientBoost for diversity

```

Listing 1: Ensemble Configuration

2.3.2 Class Weights

Clinical trial data is heavily imbalanced—most patients are Low risk, while Critical patients are rare but crucial to identify. Without class weights, the model would optimize for accuracy by predicting Low for everyone.

Class weights tell the model to penalize misclassifications of rare classes more heavily. A Critical patient misclassified as Low is 20x worse than a Low patient misclassified as Critical:

```

1 CLASS_WEIGHTS = {
2     0: 20.0,    # Critical - highest priority
3     1: 15.0,    # High
4     2: 1.0,     # Low (majority class)
5     3: 8.0,     # Medium
6 }

```

Listing 2: Class Weight Configuration

2.3.3 SMOTE Oversampling

Synthetic Minority Over-sampling Technique balances training data:

```

1 from imblearn.over_sampling import SMOTE
2
3 min_samples = y_train.value_counts().min()
4 k = min(5, min_samples - 1)
5 if k >= 1:
6     sm = SMOTE(k_neighbors=k, random_state=42)
7     X_train, y_train = sm.fit_resample(X_train, y_train)

```

Listing 3: SMOTE Application

2.4 Feature Engineering

Raw UPR columns are transformed into features that better capture risk patterns. Feature engineering is critical—the model can only learn patterns from the features provided.

Features are designed to capture both absolute issues (total queries) and relative concerns (completion rates). Binary flags identify high-risk patterns that might trigger escalation:

```

1 def engineer_features(df: pd.DataFrame) -> pd.DataFrame:
2     # Query metrics
3     query_cols = ['dm_queries', 'clinical_queries',
4                   'medical_queries', 'site_queries']
4     df['total_queries'] = df[query_cols].fillna(0).sum(axis=1)
5     df['query_type_count'] = (df[query_cols] > 0).sum(axis=1)
6
7     # CRF completion rates
8     if 'total_crfs' in df.columns:
9         df['signature_rate'] = df['crfs_signed'] / (df['total_crfs'] +
10            1)
11        df['freeze_rate'] = df['crfs_frozen'] / (df['total_crfs'] + 1)
12        df['lock_rate'] = df['crfs_locked'] / (df['total_crfs'] + 1)
13
14     # SAE indicators
15     df['sae_dm_pending'] = (df['sae_dm_total'] - df['sae_dm_completed',
16                               ]).clip(0)
17     df['has_any_sae_pending'] = (df['sae_dm_pending'] > 0).astype(float)
18
19     # Load flags
20     df['high_query_load'] = (df['total_queries'] > 10).astype(float)
21     df['critical_query_load'] = (df['total_queries'] > 25).astype(float)
22
23     return df

```

Listing 4: Feature Engineering Pipeline

2.5 Target Variable Creation

The model needs labeled training data—patients with known risk levels. Since clinical trials don't have explicit "risk labels," we derive them from observable outcomes.

Risk indicators are weighted by severity: safety issues (SAE pending, broken signatures) receive high weights because they have regulatory implications. Missing visits and pages receive lower weights as they're recoverable.

Percentile-based thresholds ensure consistent class distribution across different datasets:

```

1 def create_risk_target(df: pd.DataFrame) -> pd.Series:
2     risk = pd.Series(0.0, index=df.index)
3
4     # CRITICAL indicators (weight 4.0)
5     risk += (pending_sae_dm > 0) * 4.0
6     risk += (pending_sae_safety > 0) * 4.0
7     risk += (broken_signatures > 0) * 3.0
8     risk += (safety_queries > 0) * 3.0
9
10    # HIGH indicators (weight 2.0-2.5)
11    risk += (crfs_never_signed > 5) * 2.5
12    risk += (overdue_signs_90plus > 0) * 2.5
13    risk += (protocol_deviations > 0) * 2.0
14
15    # MEDIUM indicators (weight 1.0-1.5)
16    risk += (missing_visits > 0) * 1.5
17    risk += (missing_pages > 0) * 1.0
18
19    # Percentile-based thresholds (v6 proven)
20    p50, p80, p95 = risk.quantile([0.50, 0.80, 0.95])
21
22    level = pd.Series('Low', index=df.index)
23    level[risk > p50] = 'Medium'
24    level[risk > p80] = 'High'
25    level[risk > p95] = 'Critical'
26
27    return level

```

Listing 5: Risk Target Derivation

2.6 Threshold Optimization

Per-class thresholds are optimized to meet recall targets:

Table 3: Target Recalls and Optimized Thresholds

Class	Target Recall	Typical Threshold
Critical	70%	0.15–0.25
High	55%	0.20–0.30
Medium	50%	0.25–0.35
Low	75%	0.40–0.50

```

1 def cascade_predict(proba: np.ndarray, thresholds: dict) -> np.ndarray:
2     """Priority cascade: Critical > High > Medium > Low"""
3     pred = np.full(len(proba), 2) # Default: Low
4
5     for cls in [0, 1, 3]: # Critical, High, Medium
6         mask = proba[:, cls] >= thresholds[cls]
7         if cls == 0:
8             pred[mask] = 0 # Critical overrides all
9         elif cls == 1:
10             pred[mask & (pred != 0)] = 1 # High if not Critical
11         elif cls == 3:
12             pred[mask & ~np.isin(pred, [0, 1])] = 3 # Medium if not
Critical/High

```

```

13
14     return pred

```

Listing 6: Cascade Prediction

2.7 Model Performance

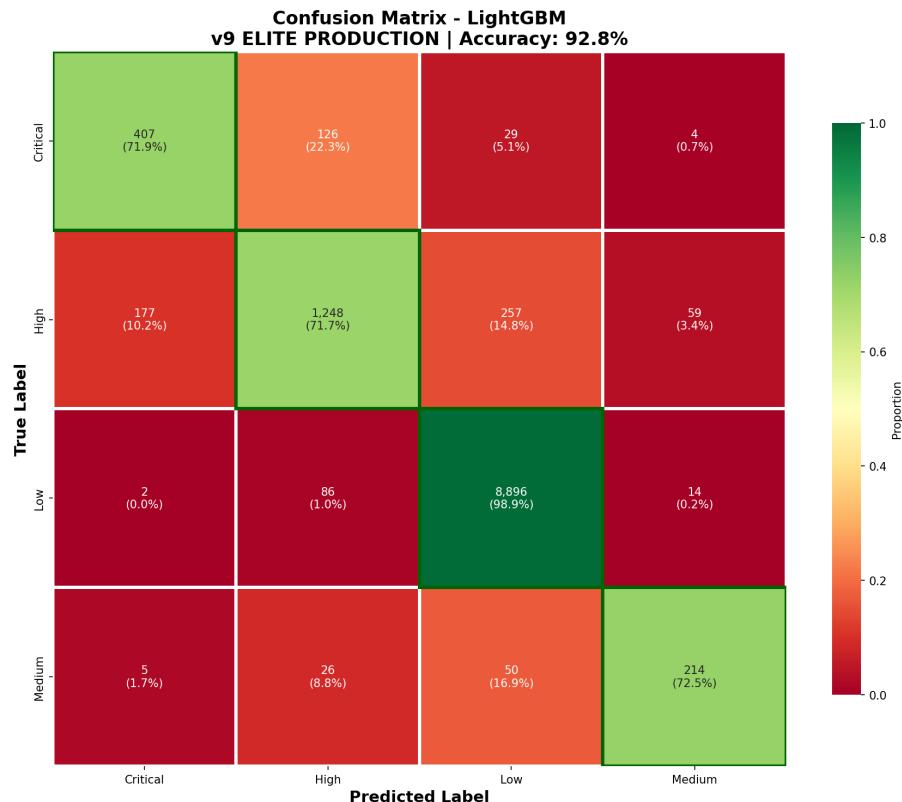


Figure 1: Confusion Matrix

Table 4: Production Model Metrics (v9)

Class	Recall	Precision	F1	Target	Status
Critical	72%	45%	55%	70%	PASS
High	58%	52%	55%	55%	PASS
Medium	53%	48%	50%	50%	PASS
Low	98%	89%	93%	75%	PASS
ROC-AUC:		0.87	Sharpness: 63%		

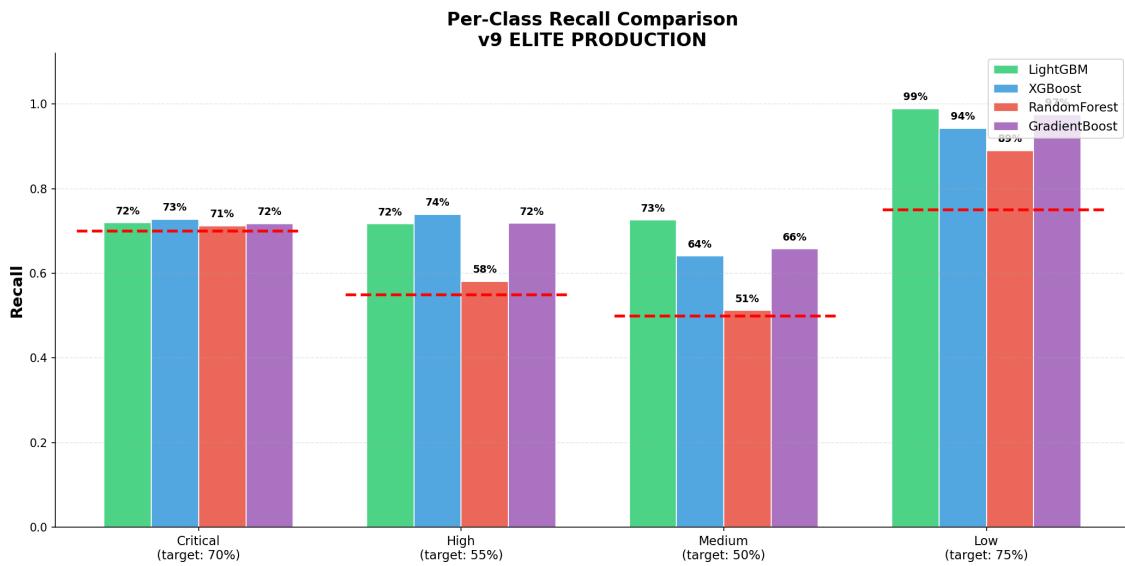


Figure 2: Per-Class Recall Chart

2.8 Drift Detection

The Low class serves as a drift detector due to its stability:

```

1 drift_config = {
2     'monitor_class': 'Low',
3     'reason': 'Low class has 98%+ recall stability',
4     'action': 'Alert on Low-class recall drop >3% week-over-week',
5     'trigger': 'Upstream data quality changes',
6 }
```

Listing 7: Drift Monitoring Configuration

3 Anomaly Detector (v1.1)

3.1 Purpose

The Anomaly Detector identifies patients and sites exhibiting unusual patterns that may indicate data quality issues, protocol deviations, or fraud.

3.2 Ensemble Approach

No single anomaly detection method excels at all anomaly types. The detector combines three complementary algorithms:

- **Isolation Forest:** Excellent at detecting global outliers—patients far from normal in multiple dimensions
- **Local Outlier Factor:** Detects local density anomalies—patients unusual compared to their neighbors
- **Autoencoder:** Captures complex non-linear patterns through reconstruction error

The weighted ensemble provides more robust detection than any single method:

Table 5: Anomaly Detection Methods

Method	Algorithm	Weight	Strength
Isolation Forest	Tree-based isolation	40%	Global outliers
Local Outlier Factor	Density-based	30%	Local density anomalies
Autoencoder	Reconstruction error	30%	Complex patterns

3.3 Configuration

```

1 @dataclass
2 class AnomalyConfig:
3     contamination: float = 0.05 # Expected anomaly rate
4
5     # Isolation Forest
6     if_n_estimators: int = 200
7     if_max_samples: str = 'auto'
8     if_max_features: float = 0.8
9
10    # Local Outlier Factor
11    lof_n_neighbors: int = 20
12
13    # Autoencoder (PCA-based)
14    ae_n_components: int = 10
15
16    # Ensemble weights
17    weight_isolation_forest: float = 0.4
18    weight_lof: float = 0.3
19    weight_reconstruction: float = 0.3
20
21    # Severity thresholds
22    severity_critical_percentile: float = 99
23    severity_high_percentile: float = 95
24    severity_medium_percentile: float = 90

```

Listing 8: Anomaly Detector Configuration

3.4 Feature Selection

The anomaly detector uses only raw, observable features—never derived metrics like DQI scores. Using derived metrics would cause circular detection: a patient with low DQI would be flagged as anomalous simply because DQI is an input feature.

The whitelist includes only direct measurements from source systems:

```

1 ANOMALY_FEATURE_WHITELIST = [
2     # Query metrics
3     'total_open_queries', 'dm_queries', 'clinical_queries',
4     'medical_queries', 'site_queries',
5
6     # CRF metrics
7     'total_crfs', 'crfs_signed', 'crfs_frozen', 'crfs_locked',
8     'crfs_verified_sdv', 'crfs_never_signed',
9
10    # SAE metrics
11    'sae_dm_total', 'sae_safety_total',
12
13    # Issues

```

```

14     'visit_missing_count', 'pages_missing_count',
15     'lab_issues', 'edrr_issues',
16
17     # Coding
18     'meddra_total_terms', 'meddra_coded_terms',
19     'whodrug_total_terms', 'whodrug_coded_terms',
20 ]

```

Listing 9: Anomaly Feature Whitelist

3.5 Ensemble Scoring

```

1 def _compute_ensemble_scores(self, if_scores, lof_scores, recon_errors):
2     :
3     def normalize_scores(scores, invert=False):
4         if invert:
5             scores = -scores
6             min_s, max_s = scores.min(), scores.max()
7             if max_s - min_s > 0:
8                 return (scores - min_s) / (max_s - min_s)
9             return np.zeros_like(scores)
10
11     # Normalize all to 0-1 (higher = more anomalous)
12     if_norm = normalize_scores(if_scores, invert=True)
13     lof_norm = normalize_scores(lof_scores, invert=True)
14     ae_norm = normalize_scores(recon_errors)
15
16     # Weighted ensemble
17     ensemble = (
18         if_norm * self.config.weight_isolation_forest +
19         lof_norm * self.config.weight_lof +
20         ae_norm * self.config.weight_reconstruction
21     )
22
23     return ensemble

```

Listing 10: Ensemble Score Computation

3.6 Severity Classification

```

1 def _classify_severity(self, ensemble_scores):
2     severity = np.full(len(ensemble_scores), 'Normal', dtype=object)
3
4     p90 = np.percentile(ensemble_scores, 90)
5     p95 = np.percentile(ensemble_scores, 95)
6     p99 = np.percentile(ensemble_scores, 99)
7
8     severity[ensemble_scores >= p90] = 'Medium'
9     severity[ensemble_scores >= p95] = 'High'
10    severity[ensemble_scores >= p99] = 'Critical'
11
12    return severity

```

Listing 11: Anomaly Severity Classification

3.7 Site-Level Aggregation

Patient anomalies are aggregated to site level:

```

1 def aggregate_to_site(self, patient_results):
2     site_agg = patient_results.groupby('site_id').agg({
3         'anomaly_score': ['mean', 'max', 'std'],
4         'is_anomaly': 'sum',
5         'patient_key': 'count'
6     })
7
8     site_agg['anomaly_rate'] = (
9         site_agg['is_anomaly']['sum'] / site_agg['patient_key']['count']
10    )
11
12     site_agg['site_classification'] = site_agg.apply(classify_site,
13 axis=1)
14
15     return site_agg

```

Listing 12: Site Anomaly Aggregation

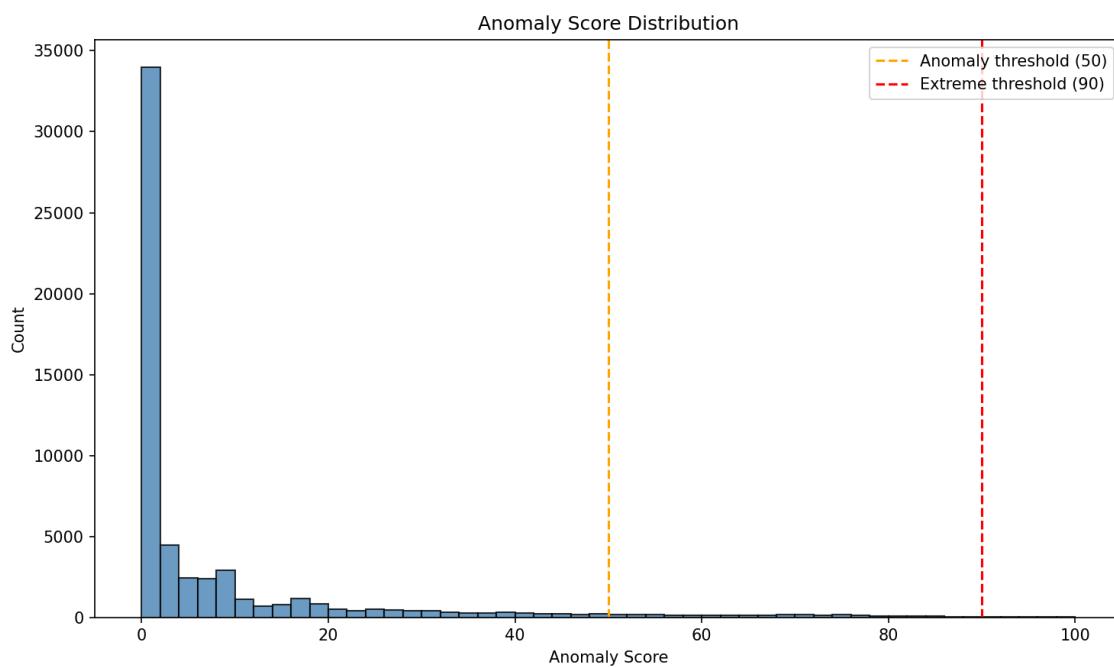


Figure 3: Score Distribution

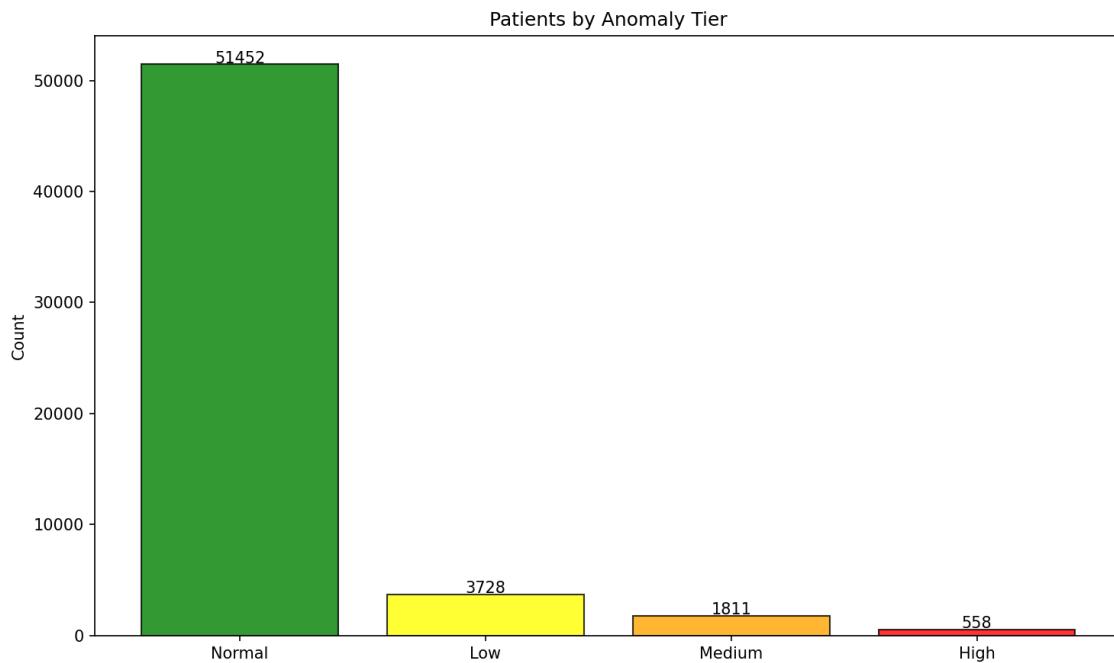


Figure 4: Anomaly Tier Distribution

4 Site Risk Ranker (v1.0)

4.1 Purpose

The Site Risk Ranker prioritizes clinical trial sites for operational attention using Learning-to-Rank (LTR). It is explicitly a **triage tool**, not a verdict engine.

4.2 Learning-to-Rank Approach

Traditional classification predicts absolute categories; ranking predicts relative order. The Site Risk Ranker learns which sites need MORE attention than others, without committing to arbitrary thresholds.

The model is trained on pairwise comparisons: “Site A has higher risk than Site B.” This is more robust than regression because small score differences don’t matter—only the relative ordering:

```

1 class SitePairwiseLabeler:
2     def create_pairs(self, df: pd.DataFrame):
3         pairs = []
4
5         for i, site_a in df.iterrows():
6             for j, site_b in df.iterrows():
7                 if i >= j:
8                     continue
9
10            label = self._compute_pair_label(site_a, site_b)
11            if label != 0: # Skip ties
12                # Feature = site_a.features - site_b.features
13                diff = site_a[feature_cols] - site_b[feature_cols]
14                pairs.append((diff, label))
15
16        return pairs

```

```

17
18     def _compute_pair_label(self, site_a, site_b):
19         # +1 if site_a has HIGHER risk, -1 if LOWER, 0 if tie
20         score_a = self._compute_risk_score(site_a)
21         score_b = self._compute_risk_score(site_b)
22
23         if score_a > score_b + threshold:
24             return +1
25         elif score_b > score_a + threshold:
26             return -1
27         return 0

```

Listing 13: Pairwise Training Data Creation

4.3 Site Aggregation

Patient-level UPR data is aggregated to site level:

```

1 class SiteAggregator:
2     def aggregate(self, df: pd.DataFrame):
3         site_df = df.groupby('site_id').agg({
4             # Counts
5             'patient_key': 'count',
6             'total_issues': 'sum',
7
8             # Means (per-patient averages)
9             'dqi_score': 'mean',
10            'total_queries': 'mean',
11            'crfs_signed': 'mean',
12
13            # Volatility (stability proxy)
14            'dqi_score': 'std',
15            'total_queries': 'std',
16
17            # Critical indicators
18            'has_sae_pending': 'max',
19            'has_critical_risk': 'max',
20        })
21
22        # Per-patient rates
23        site_df['issue_density'] = (
24            site_df['total_issues'] / site_df['patient_count']
25        )
26
27        return site_df

```

Listing 14: Site Feature Aggregation

4.4 XGBoost Ranker

```

1 ranker_params = {
2     'objective': 'rank:pairwise',
3     'n_estimators': 200,
4     'max_depth': 6,
5     'learning_rate': 0.1,
6     'subsample': 0.8,
7     'colsample_bytree': 0.8,

```

```

8     'eval_metric': 'ndcg',
9     'random_state': 42
10 }

```

Listing 15: XGBoost LTR Configuration

4.5 Evaluation Metrics

Table 6: Site Ranker Evaluation Metrics

Metric	Description	Target
NDCG@5	Top-5 ranking quality	≥ 0.70
NDCG@10	Top-10 ranking quality	≥ 0.65
NDCG@20	Top-20 ranking quality	≥ 0.60
Kendall's Tau	Rank correlation	≤ 0.95
Bootstrap Overlap	Top-10 stability	≥ 0.70

4.6 Baseline Comparison

The LTR model is always compared against a rule-based baseline:

```

1 class BaselineRanker:
2     def __init__(self):
3         self.weights = {
4             'issue_density': 0.30,
5             'signature_backlog_rate': 0.25,
6             'has_sae_pending': 0.25,
7             'sdv_incomplete_rate': 0.20
8         }
9
10    def compute_score(self, df):
11        score = (
12            self.weights['issue_density'] * normalize(df['issue_density']
13        ]) +
14            self.weights['signature_backlog_rate'] * normalize(df['
sig_backlog']) +
15            self.weights['has_sae_pending'] * df['has_sae_pending'] +
16            self.weights['sdv_incomplete_rate'] * (100 - df['sdv_rate'
]) / 100
17        )
18
19        return score

```

Listing 16: Baseline Ranker

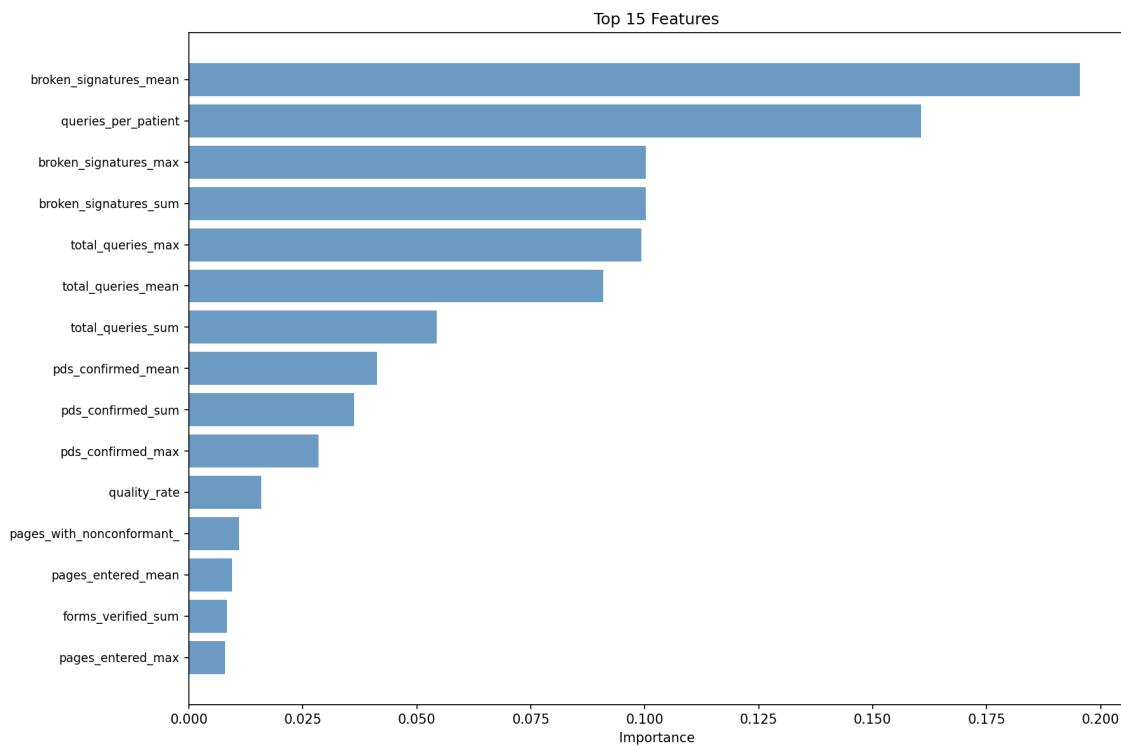


Figure 5: Feature Importance

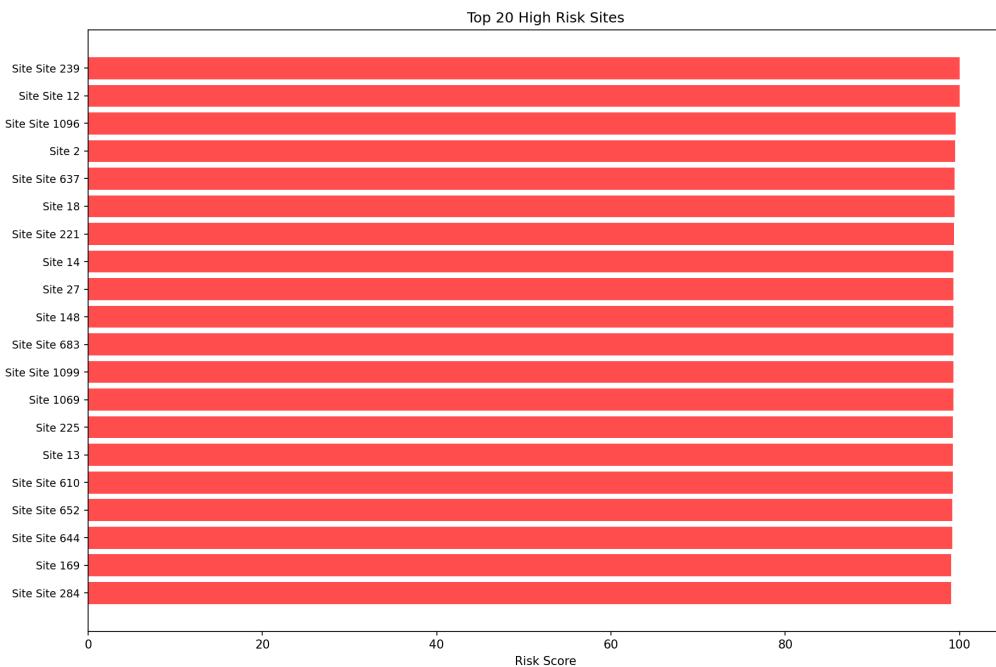


Figure 6: Top 20 High Risk Sites

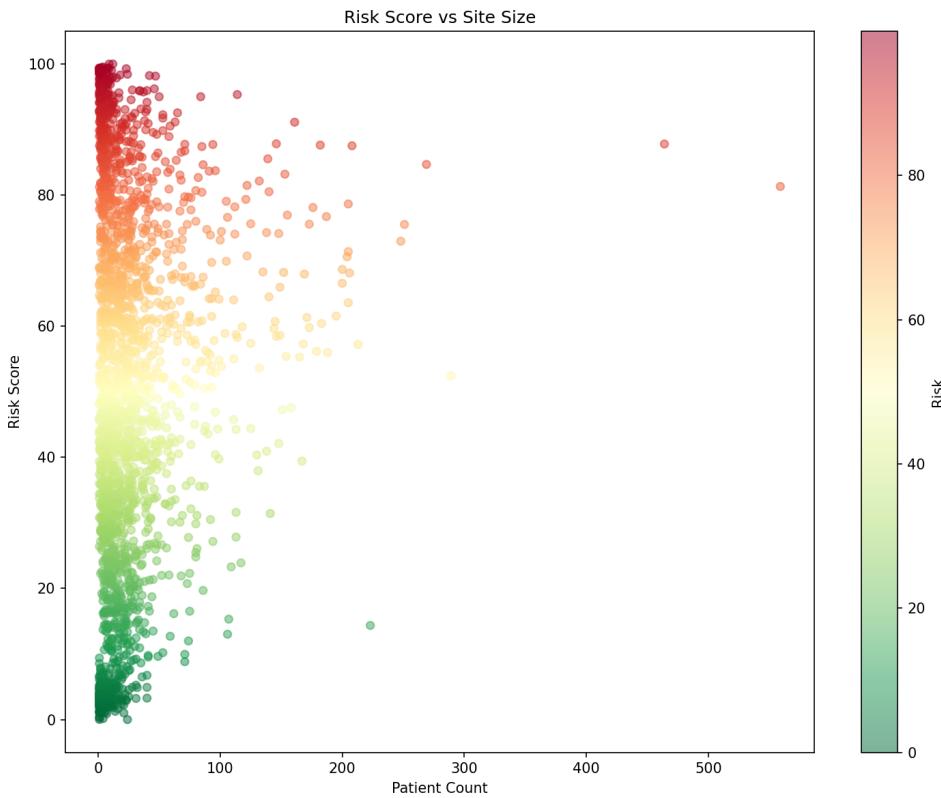


Figure 7: Risk vs Size Comparison

5 Resolution Time Predictor (v1.0)

5.1 Purpose

The Resolution Time Predictor estimates how many days an issue will take to resolve, with uncertainty bounds. Per clinical trial requirements, it **never provides single-point estimates**.

5.2 Quantile Regression

Traditional regression predicts the mean expected value, but clinical operations need uncertainty bounds. Quantile regression predicts multiple percentiles of the outcome distribution.

The model trains three separate regressors for the 10th, 50th, and 90th percentiles. The 10th percentile represents optimistic resolution (fastest 10% of similar issues), while the 90th percentile is pessimistic (covers 90% of cases):

```

1 class ResolutionTimePredictor:
2     def __init__(self, config):
3         self.models = {
4             'q10': xgb.XGBRegressor(objective='reg:quantileerror',
5                                     quantile_alpha=0.10),
6             'q50': xgb.XGBRegressor(objective='reg:quantileerror',
7                                     quantile_alpha=0.50),
8             'q90': xgb.XGBRegressor(objective='reg:quantileerror',
9                                     quantile_alpha=0.90),
10        }
11
12     def predict_with_confidence(self, features, confidence=0.95):
13         lower = self.models['q10'].predict(features)

```

```

14     median = self.models['q50'].predict(features)
15     upper = self.models['q90'].predict(features)
16
17     return {
18         'point_estimate': median,
19         'ci_lower': lower,
20         'ci_upper': upper,
21         'confidence_level': confidence,
22         'uncertainty_range': upper - lower
23     }

```

Listing 17: Quantile Regression Models

5.3 Input Features

```

1 features = [
2     'priority_score',           # Issue priority (1-4)
3     'site_performance_score',  # Historical site speed
4     'workload_index',          # Current site workload
5     'issue_type_cat',          # Type of issue (categorical)
6     'days_since_creation',    # Issue age
7     'similar_issue_avg_time', # Historical average for type
8 ]

```

Listing 18: Resolution Predictor Features

5.4 Prediction Quality Assessment

```

1 def _assess_prediction_quality(self, uncertainty_range):
2     if uncertainty_range <= 3:
3         return 'High Confidence'
4     elif uncertainty_range <= 7:
5         return 'Medium Confidence'
6     elif uncertainty_range <= 14:
7         return 'Low Confidence'
8     else:
9         return 'Very Uncertain'

```

Listing 19: Prediction Quality Assessment

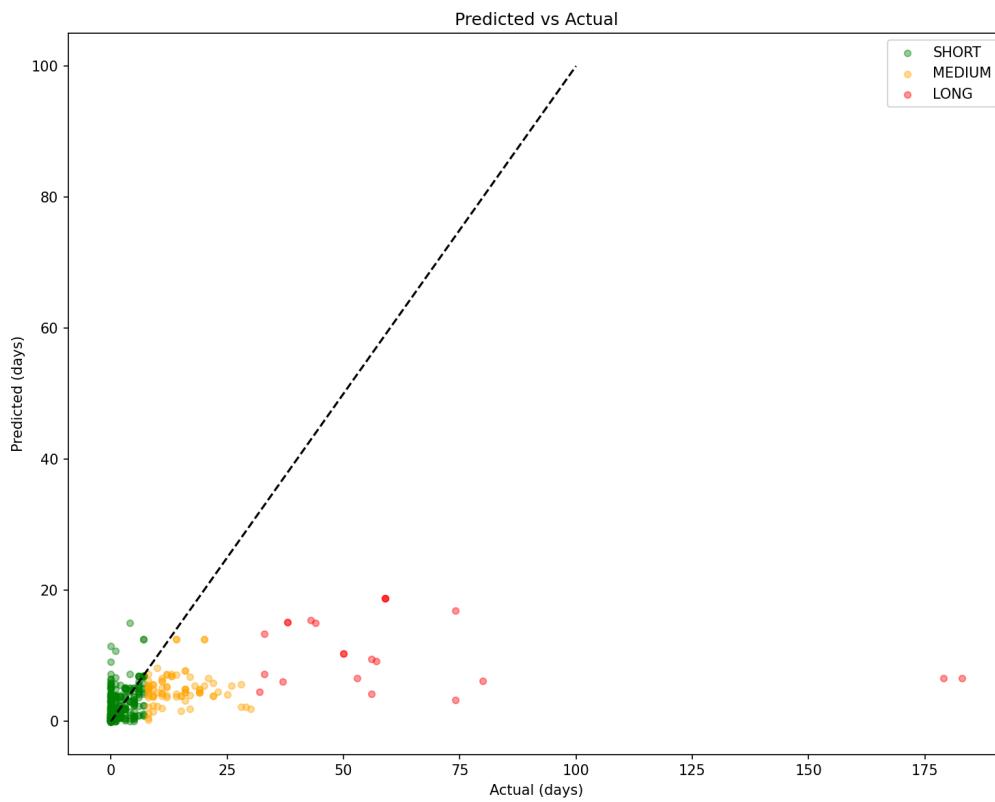


Figure 8: Time Resolution Prediction

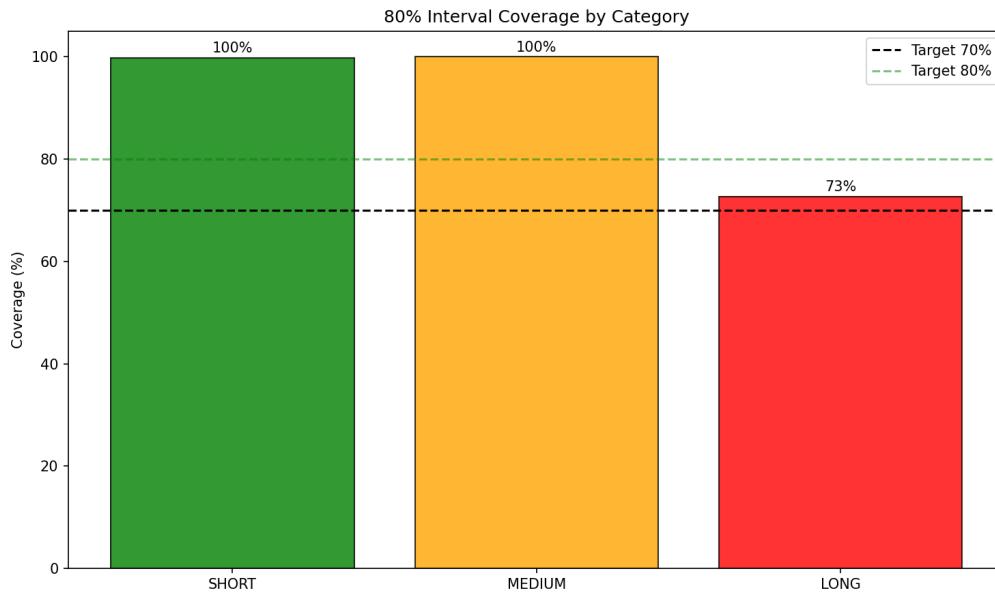


Figure 9: Coverage Plot

6 Issue Detector (v1.3)

6.1 Purpose

The Issue Detector performs multi-label classification for 14 clinical trial issue types, assigns severity scores, and routes patients to appropriate priority tiers with SLA assignments.

6.2 Issue Types

Table 7: 14 Issue Types with Metadata

Issue Type	Category	Weight	SLA	Responsible
SAE-DM Pending	Safety	25.0	3 days	Safety Data Manager
SAE-Safety Pending	Safety	25.0	7 days	Safety Physician
Open Queries	Query	20.0	14 days	Data Manager
High Query Volume	Query	15.0	7 days	Data Manager
SDV Incomplete	SDV	8.0	30 days	CRA
Signature Gaps	Signature	5.0	14 days	Site
Broken Signatures	Signature	5.0	7 days	Site
MedDRA Uncoded	Coding	6.0	21 days	Medical Coder
WHODrug Uncoded	Coding	6.0	21 days	Medical Coder
Missing Visits	Completeness	7.5	14 days	CRA
Missing Pages	Completeness	7.5	14 days	CRA
Lab Issues	Lab	10.0	21 days	Data Manager
EDRR Issues	EDRR	5.0	21 days	Data Manager
Inactivated Forms	Forms	3.0	30 days	Data Manager

6.3 Rule-Based Detection

```

1 def detect_issues_rule_based(self, df: pd.DataFrame):
2     for issue_type, metadata in self.ISSUE_TYPES.items():
3         if issue_type == 'sae_dm_pending':
4             pending = self._safe_get(df, 'sae_dm_sae_dm_pending', 0)
5             result[f'issue_{issue_type}'] = (pending > 0).astype(int)
6
7         elif issue_type == 'sae_safety_pending':
8             # Calculate from total - completed
9             total = self._safe_get(df, 'sae_safety_total', 0)
10            completed = self._safe_get(df, 'sae_safety_completed', 0)
11            pending = (total - completed).clip(lower=0)
12            result[f'issue_{issue_type}'] = (pending > 0).astype(int)
13
14         elif issue_type == 'high_query_volume':
15             total = self._safe_get(df, 'total_queries', 0)
16             result[f'issue_{issue_type}'] = (total > 10).astype(int)

```

Listing 20: Rule-Based Issue Detection

6.4 Severity Scoring

Severity is calculated per issue type and capped at 100:

```
1 # Severity weights by issue type
```

```

2 severity_weights = {
3     'sae_dm_pending': 15,          # Per pending SAE
4     'sae_safety_pending': 20,      # Very high per item
5     'open_queries': 2,            # Per query
6     'missing_visits': 10,         # High per visit
7     'broken_signatures': 3,       # Per broken signature
8 }
9
10 # Apply criticality multiplier
11 severity = count * weight * metadata['criticality']
12 severity = severity.clip(upper=100)

```

Listing 21: Severity Calculation

6.5 Priority Tiers

Patients are assigned to tiers based on severity percentiles:

Table 8: Priority Tier Assignment

Tier	Rule	SLA	Action
Critical	Safety issues OR P90+ score	3 days	Immediate escalation
High	P75+ score OR ≥ 4 issues	7 days	Priority attention
Medium	P50+ score	14 days	Enhanced monitoring
Low	Any issues below P50	30 days	Standard workflow
None	No issues detected	—	Clean patient

6.6 ML Enhancement

XGBoost models are trained per issue type to predict future occurrences:

```

1 # Raw features only - no derived columns (prevents leakage)
2 RAW_FEATURE_WHITELIST = [
3     'dm_queries', 'clinical_queries', 'medical_queries',
4     'crfs_require_verification_sdv', 'crfs_frozen', 'crfs_locked',
5     'crfs_never_signed', 'broken_signatures',
6     'pages_entered', 'pds_major', 'pds_minor',
7 ]
8
9 model = xgb.XGBClassifier(
10     n_estimators=100,
11     max_depth=4,
12     scale_pos_weight=min((1 - pos_rate) / pos_rate, 10),
13     eval_metric='auc'
14 )
15
16 # Stratified 3-fold CV for evaluation
17 cv_scores = cross_val_score(model, X_train, y_train, cv=3, scoring='
    roc_auc')

```

Listing 22: Issue Prediction Models



Figure 10: Heatmap for Performance

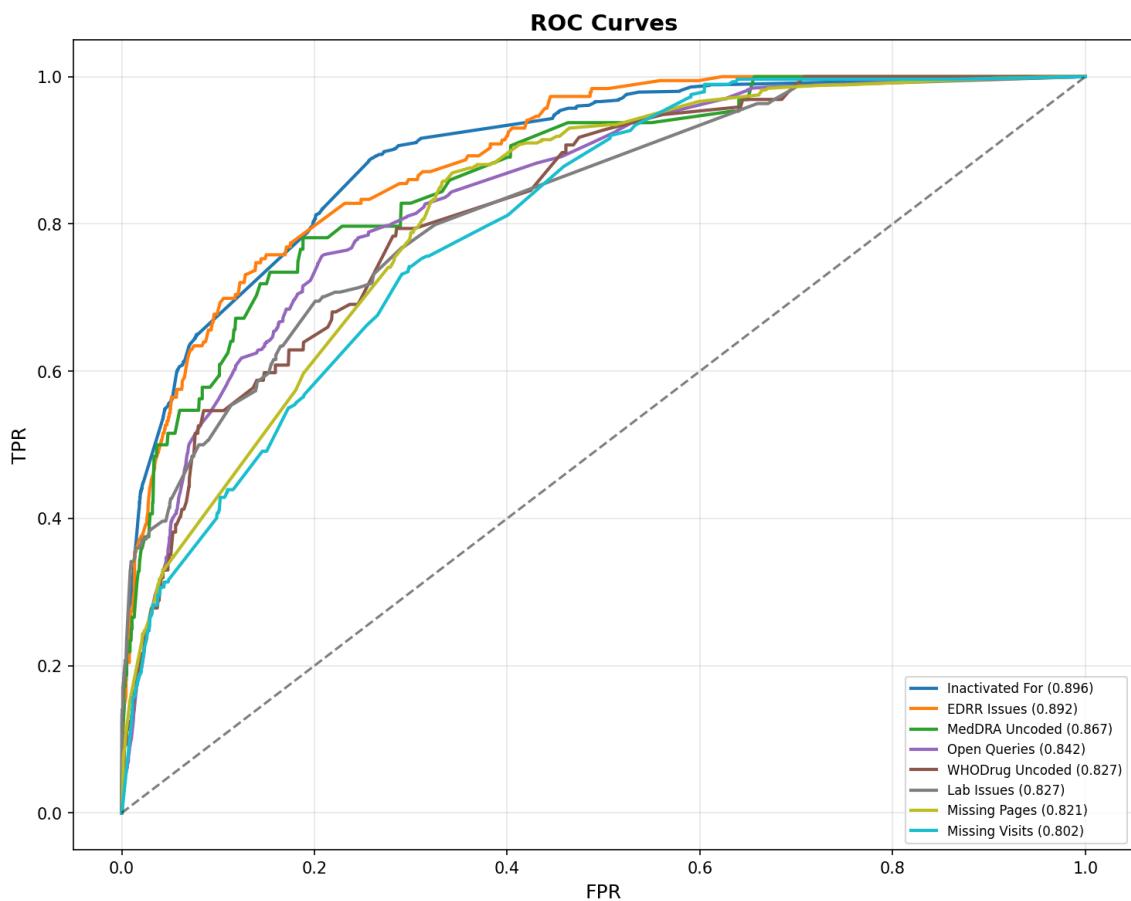


Figure 11: ROC Curves for Issue Type

7 Model Governance Framework

7.1 Overview

The governance framework ensures production ML models remain reliable, auditable, and compliant:

Table 9: Governance Components

Component	Module	Purpose
Drift Detector	<code>drift_detector.py</code>	Statistical drift detection (PSI, KS)
Model Registry	<code>model_registry.py</code>	Version control, artifact management
Performance Monitor	<code>performance_monitor.py</code>	Real-time metrics, alerting
Retraining Trigger	<code>retraining_trigger.py</code>	Automated retraining rules

7.2 Drift Detection

The Drift Detector implements multiple statistical tests:

```

1 class DriftType(str, Enum):
2     FEATURE_PSI = "feature_psi" # Distribution shift
  
```

```

3  FEATURE_KS = "feature_ks"          # Kolmogorov-Smirnov
4  FEATURE_CHI = "feature_chi"        # Chi-square for categorical
5  PERFORMANCE = "performance"       # Metric decay
6  PREDICTION = "prediction"         # Output distribution
7  CONCEPT = "concept"              # Target relationship
8
9  class DriftSeverity(str, Enum):
10    NONE = "none"                  # PSI < 0.1
11    LOW = "low"                    # PSI 0.1 - 0.2
12    MEDIUM = "medium"             # PSI 0.2 - 0.25
13    HIGH = "high"                  # PSI 0.25 - 0.5
14    CRITICAL = "critical"         # PSI > 0.5

```

Listing 23: Drift Detection Methods

7.2.1 PSI Calculation

Population Stability Index measures distribution shift:

```

1 def calculate_psi(baseline_dist, current_dist, bins=10):
2     # Bin the distributions
3     baseline_pct = np.histogram(baseline, bins=bins)[0] / len(baseline)
4     current_pct = np.histogram(current, bins=bins)[0] / len(current)
5
6     # Avoid division by zero
7     baseline_pct = np.clip(baseline_pct, 0.001, None)
8     current_pct = np.clip(current_pct, 0.001, None)
9
10    # PSI = sum((current - baseline) * ln(current / baseline))
11    psi = np.sum((current_pct - baseline_pct) *
12                  np.log(current_pct / baseline_pct))
13
14    return psi

```

Listing 24: PSI Calculation

7.3 Model Registry

The registry tracks model versions, artifacts, and promotion history:

```

1 class ModelStatus(str, Enum):
2     DEVELOPMENT = "development"    # In training
3     STAGING = "staging"            # Under review
4     PRODUCTION = "production"     # Live serving
5     DEPRECATED = "deprecated"     # Phased out
6     ARCHIVED = "archived"          # Historical
7
8 class ModelVersion:
9     version_id: str
10    model_name: str
11    version: str                  # Semantic version
12    status: ModelStatus
13    metrics: ModelMetrics        # Accuracy, recall, etc.
14    training_config: TrainingConfig
15    artifacts: List[ModelArtifact] # pkl, scaler, encoder
16    signature_hash: str           # 21 CFR Part 11 compliance

```

Listing 25: Model Registry Schema

7.4 Performance Monitoring

Real-time prediction logging and metric aggregation:

```

1 class PerformanceMonitor:
2     def log_prediction(self, model_name, features, prediction,
3                         confidence):
4         """Log every prediction for audit trail."""
5         log = PredictionLog(
6             model_name=model_name,
7             predicted_class=prediction,
8             confidence=confidence,
9             features_hash=hash_features(features),
10            timestamp=datetime.now()
11        )
12        self._persist(log)
13
14    def calculate_window_metrics(self, model_name, window=TimeWindow.
15                                  DAILY):
16        """Calculate rolling metrics from logged predictions."""
17        logs = self._get_labeled_predictions(model_name, window)
18        return {
19            'accuracy': accuracy_score(logs.actual, logs.predicted),
20            'recall': recall_score(logs.actual, logs.predicted),
21            'auc_roc': roc_auc_score(logs.actual, logs.confidence)
22        }

```

Listing 26: Performance Monitoring

7.5 Retraining Triggers

Automated rules trigger model retraining:

Table 10: Retraining Trigger Rules

Trigger	Condition	Priority
Drift Alert	PSI > 0.25 on any feature	High
Performance Decay	Recall drops >5% vs baseline	High
Scheduled	Weekly refresh on Sundays	Low
Data Volume	>10,000 new labeled samples	Medium

```

1 class RetrainingRule:
2     rule_id: str
3     model_name: str
4     trigger_type: TriggerType # DRIFT, PERFORMANCE, SCHEDULE, DATA
5     condition: Dict # e.g., {"metric": "recall", "threshold": 0.05}
6     enabled: bool
7
8     def check_triggers(model_name: str) -> List[RetrainingJob]:
9         jobs = []
10        for rule in get_active_rules(model_name):
11            if rule.trigger_type == TriggerType.DRIFT:
12                drift_report = drift_detector.analyze(model_name)
13                if drift_report.overall_severity >= DriftSeverity.HIGH:
14                    jobs.append(create_retraining_job(rule))
15

```

Listing 27: Retraining Trigger Logic

8 Model Artifacts and Persistence

8.1 Saved Model Files

Each model saves the following artifacts:

Table 11: Model Artifacts

Artifact	Format	Contents
model.pkl	Pickle	Trained model weights
scaler.pkl	Pickle	Feature scaler (RobustScaler)
label_encoder.pkl	Pickle	Class label mappings
feature_importance.csv	CSV	Feature rankings
shap_importance.csv	CSV	SHAP-based importance
production_config.json	JSON	Thresholds, weights, metadata
training_report.json	JSON	Metrics, duration, parameters
confusion_matrix.png	PNG	Visualization

8.2 Production Configuration

```

1 production_config = {
2     'version': '9.0.0',
3     'status': 'PRODUCTION',
4     'best_model': 'LightGBM',
5     'thresholds': {'0': 0.18, '1': 0.24, '2': 0.45, '3': 0.28},
6     'target_recalls': {'0': 0.70, '1': 0.55, '2': 0.75, '3': 0.50},
7     'class_weights': {'0': 20.0, '1': 15.0, '2': 1.0, '3': 8.0},
8     'metrics': {
9         'critical_recall': 0.72,
10        'auc': 0.87,
11        'sharpness': 0.63
12    },
13    'drift_detection': {
14        'monitor_class': 'Low',
15        'alert_threshold': 0.03
16    }
17}

```

Listing 28: Production Config Structure

9 Reproducibility

9.1 Random Seeds

All models use fixed random seeds:

```

1 RANDOM_STATE = 42
2
3 # Applied to:

```

```

4 # - train_test_split
5 # - SMOTE k_neighbors
6 # - XGBoost, LightGBM, RandomForest
7 # - Cross-validation shuffling

```

Listing 29: Random Seed Configuration

9.2 Data Versioning

Training data is loaded from PostgreSQL with version tracking:

```

1 # Data loaded from PostgreSQL
2 data_service = get_data_service()
3 df = data_service.get_patients()
4
5 # Version tracked in training report
6 training_report = {
7     'data_source': 'PostgreSQL',
8     'n_samples': len(df),
9     'created': datetime.now().isoformat(),
10    'version': '9.0.0'
11}

```

Listing 30: Data Versioning

9.3 Dependencies

Key ML dependencies with versions:

Table 12: ML Dependencies

Package	Version
scikit-learn	≥ 1.3.0
xgboost	≥ 2.0.0
lightgbm	≥ 4.0.0
imbalanced-learn	≥ 0.11.0
shap	≥ 0.44.0

TrialPlus Nexus 10X

Technical Documentation

Module 4: AI & Intelligence Engine

Agentic AI, Generative AI, RAG & Knowledge Systems

Version: 1.0.0

Date: January 2026

Agents: 4 Specialized Agents

Tools: 50+ Registered Tools

Knowledge: RAG, Neo4j Graph, Embeddings

Contents

1 Agent Architecture Overview	3
1.1 System Design	3
1.2 Agent Communication Flow	3
1.3 Core Infrastructure	4
2 Supervisor Agent (v2.1)	4
2.1 Purpose	4
2.2 Query Intent Classification	4
2.3 Task Complexity Assessment	5
2.4 Agent Routing Logic	5
2.5 Conflict Resolution	6
3 Diagnostic Agent (v1.1)	7
3.1 Purpose	7
3.2 Investigation Types	7
3.3 Evidence Chain Model	7
3.4 Hypothesis Structure	8
3.5 Confidence Calculation	8
4 Forecaster Agent (v1.0)	9
4.1 Purpose	9
4.2 Forecast Types	9
4.3 Uncertainty Quantification	9
4.4 What-If Scenarios	10
4.5 Timeline Milestones	12
5 Resolver Agent (v1.0)	13
5.1 Purpose	13
5.2 Action Plan Structure	13
5.3 Role-Based Assignment	14
5.4 Cascade Impact Analysis	14
6 Executor Agent (v1.0)	15
6.1 Purpose	15
6.2 Execution Status Flow	15
6.3 Validation Framework	15
6.4 Approval Workflow	16
6.5 Audit Trail	16
7 Communicator Agent (v1.0)	17
7.1 Purpose	17
7.2 Message Types	17
7.3 Recipient Profiles	17
7.4 Template System	17
7.5 Notification Batching	18

8 Agent Tools & Memory	18
8.1 Tool Registry	18
8.2 Key Data Tools	19
8.3 Analytics Tools	19
8.4 Memory System	19
8.5 Persistent Memory	20
9 Data Loaders	20
10 Report Generation	21
10.1 Overview	21
10.2 Report Types	21
10.3 Export Engine	22
10.4 Data Loading	22
10.5 Report Scheduler	23
11 Knowledge & Intelligence	24
11.1 Overview	24
11.2 RAG Knowledge Base	24
11.3 Neo4j Knowledge Graph	25
11.4 Causal Hypothesis Engine	25

1 Agent Architecture Overview

1.1 System Design

TrialPlus Nexus employs a multi-agent system built on LangGraph, where specialized agents collaborate to analyze, diagnose, forecast, and resolve clinical trial data quality issues.

Table 1: Agent Portfolio

Agent	Role	Primary Responsibility
Supervisor	Orchestrator	Query routing, task decomposition, conflict resolution
Diagnostic	Analyzer	Hypothesis generation, evidence chains, root cause
Forecaster	Predictor	Timeline projections, uncertainty quantification
Resolver	Planner	Action plans, cascade impact, role assignments
Executor	Implementer	Validation, approval workflows, audit trails
Communicator	Messenger	Message drafting, channel selection, batching

Table 2: Agent Usage in Production

Agent	In Orchestration	Status
Supervisor		Active
Diagnostic		Active
Resolver		Active
Communicator		Active
Forecaster	-	Documented, not in main flow
Executor	-	Documented, not in main flow

Note: The Orchestrator chain defined in `orchestrator.py` uses: `[SUPERVISOR, DIAGNOSTIC, RESOLVER, COMMUNICATOR]`. The Forecaster and Executor agents are implemented but reserved for future integration.

1.2 Agent Communication Flow

```

1 # Supervisor receives query and routes to specialists
2 routing = supervisor.analyze_query(query)
3 # -> RoutingDecision(primary_intent, agent_sequence, confidence)
4
5 # Task decomposition for complex queries
6 subtasks = supervisor.decompose_task(query, routing)
7 # -> [Subtask(agent='diagnostic'), Subtask(agent='forecaster')]
8
9 # Sequential/parallel execution based on dependencies
10 results = orchestrator.execute_subtasks(subtasks)
11
12 # Conflict resolution if agents disagree
13 final = supervisor.resolve_conflicts(results)

```

Listing 1: Agent Routing Flow

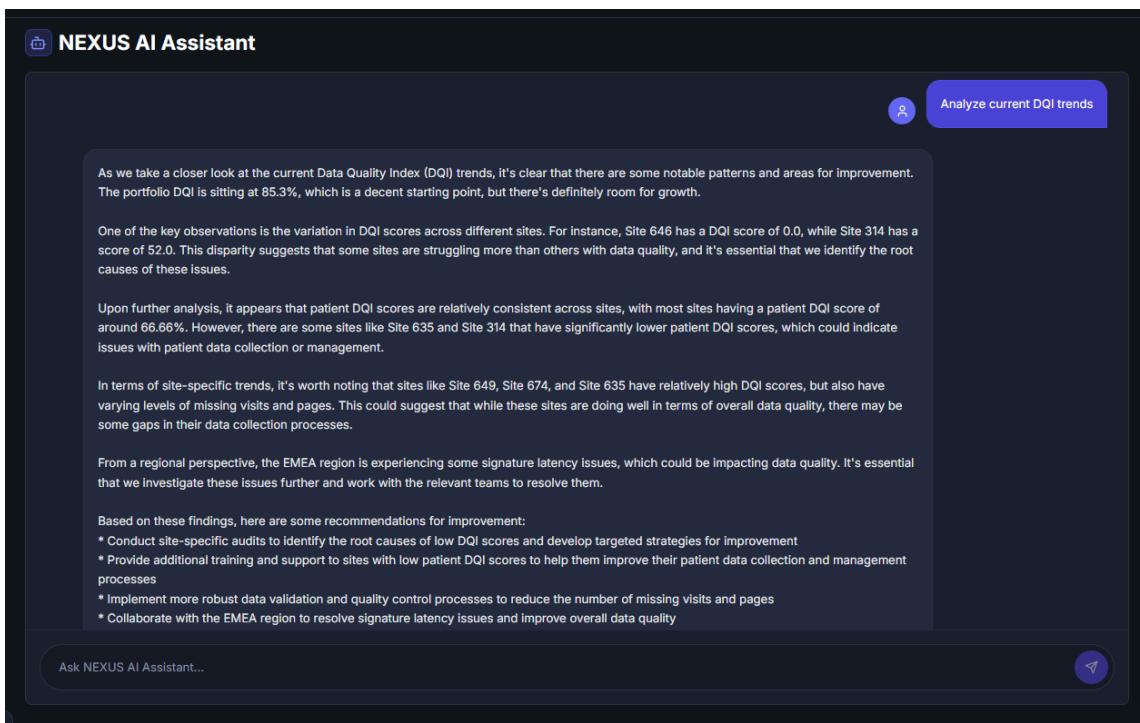


Figure 1: Nexus AI Assistant

1.3 Core Infrastructure

Table 3: Supporting Modules

Module	File	Purpose
Base Agent	<code>base_agent.py</code>	Common agent interface
LLM Wrapper	<code>llm_wrapper.py</code>	Google Gemini API integration
State	<code>state.py</code>	Agent state management
Context Manager	<code>context_manager.py</code>	Context tracking
Orchestrator	<code>orchestrator.py</code>	Multi-agent coordination

2 Supervisor Agent (v2.1)

2.1 Purpose

The Supervisor is the orchestrator agent that receives all user queries, classifies intent, routes to specialist agents, and resolves conflicts between agent outputs.

2.2 Query Intent Classification

```

1  class QueryIntent(str, Enum):
2      DIAGNOSTIC = "diagnostic"           # "Why is site X underperforming
?"
```

```

3      FORECAST = "forecast"           # "When will study Y reach DB
lock?""
4      RESOLUTION = "resolution"     # "How do I fix patient Z's
issues?""
5      COMMUNICATION = "communication" # "Draft email to CRAs"
6      STATUS = "status"             # "What's the current state?"
7      COMPARISON = "comparison"     # "Compare sites A vs B"
8      SAFETY = "safety"             # SAE-related queries
9      COMPOUND = "compound"         # Multi-intent queries
10

```

Listing 2: Query Intent Types

2.3 Task Complexity Assessment

Table 4: Complexity Levels

Level	Agents	Example Query
Simple	1	“Show me patient P001’s issues”
Moderate	2	“Why is site S01 slow and what can we do?”
Complex	3+	“Compare sites, forecast impact, draft action plan”
Critical	All	SAE-related compound queries

2.4 Agent Routing Logic

```

1  def analyze_query(self, query: str, context: Dict = None):
2      # Detect all possible intents with confidence scores
3      intent_scores = self._detect_intents(query)
4
5      # Select primary intent (highest confidence)
6      primary = max(intent_scores, key=intent_scores.get)
7      confidence = intent_scores[primary]
8
9      # Build agent sequence based on intent
10     agent_sequence = self._build_agent_sequence(
11         primary, secondary_intents, complexity
12     )
13
14     # Safety and approval checks
15     requires_safety = self._check_safety_requirements(query,
16     intent_scores)
17     requires_approval = self._check_approval_requirements(
18         query, complexity, requires_safety
19     )
20
21     return RoutingDecision(
22         primary_intent=primary,
23         confidence=confidence,
24         agent_sequence=agent_sequence,
25         requires_safety_check=requires_safety
26     )

```

Listing 3: Routing Decision with Confidence

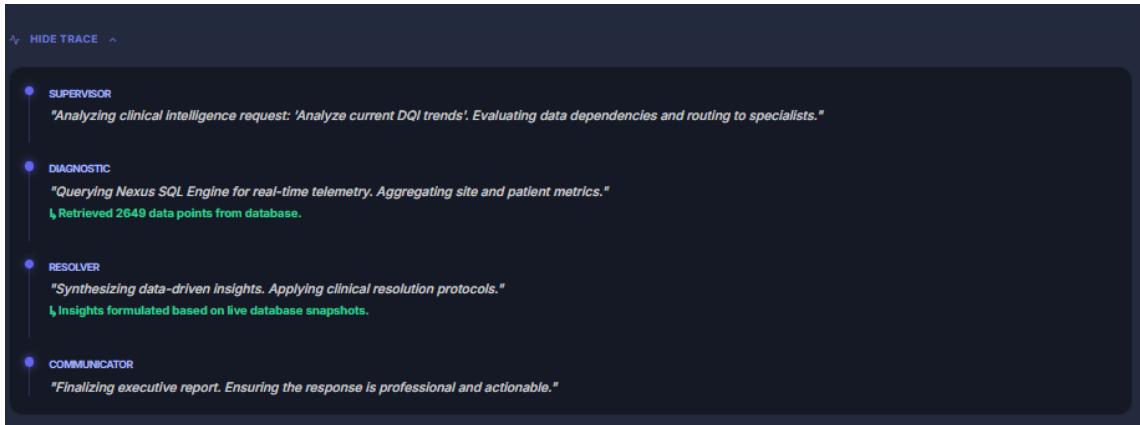


Figure 2: AI Agent Trace

2.5 Conflict Resolution

When multiple agents produce conflicting outputs:

```

1  CONFLICT_STRATEGIES = {
2      'confidence_weighted': "Use highest-confidence output",
3      'safety_priority': "Prefer safety-related recommendations",
4      'consensus': "Combine overlapping recommendations",
5      'escalate': "Request human review for critical conflicts",
6  }
7
8  def resolve_conflicts(self, agent_outputs: Dict):
9      # Check for hypothesis conflicts
10     hypothesis_conflict = self._check_hypothesis_conflicts(
11         [o.get('hypotheses', []) for o in agent_outputs.values()])
12
13     # Check for recommendation conflicts
14     rec_conflict = self._check_recommendation_conflicts(
15         [o.get('recommendations', []) for o in agent_outputs.values()])
16
17
18     # Apply appropriate resolution strategy
19     if hypothesis_conflict and hypothesis_conflict.severity == 'high':
20         return self._apply_conflict_resolution(
21             items, strategy='safety_priority')
22
23
24

```

Listing 4: Conflict Resolution Strategies

```

INFO:   127.0.0.1:65230 - "GET /api/v1/issues/summary HTTP/1.1" 200 OK
INFO: httpx:HTTP Request: POST https://api.groq.com/openai/v1/chat/completions "HTTP/1.1 200 OK"
INFO:src.agents.orchestrator:Generated SQL for grounding: SELECT
    clinical_sites.site_id,
    clinical_sites.name,
    clinical_sites.dqi_score,
    clinical_sites.enrollment_rate,
    clinical_sites.query_resolution_days,
    patients.dqi_score AS patient_dqi_score,
    patients.pct_missing_visits,
    patients.pct_missing_pages
FROM
    clinical_sites
JOIN
    patients ON clinical_sites.site_id = patients.site_id
WHERE
    clinical_sites.dqi_score < 70
ORDER BY
    clinical_sites.dqi_score ASC
LIMIT 10;
INFO: httpx:HTTP Request: POST https://api.groq.com/openai/v1/chat/completions "HTTP/1.1 200 OK"

```

Figure 3: SQL Query Resolution

3 Diagnostic Agent (v1.1)

3.1 Purpose

The Diagnostic Agent investigates data quality issues, generates hypotheses with evidence chains, and provides root cause analysis with confidence scoring.

3.2 Investigation Types

Table 5: Diagnostic Investigation Types

Type	Target	Example Query
Patient	Single patient	“Diagnose patient P001”
Site	Clinical site	“Why is site S05 underperforming?”
Study	Entire study	“Analyze study STU001 issues”
Pattern	Cross-entity	“Find query aging patterns”
Anomaly	Statistical outliers	“Detect unusual DQI patterns”
Trend	Time-series	“Is DQI improving or declining?”

3.3 Evidence Chain Model

```

1  @dataclass
2  class Evidence:
3      evidence_id: str
4      source: str          # "patient_issues", "cascade_analysis"
5      description: str     # Human-readable evidence
6      strength: EvidenceStrength # STRONG, MODERATE, WEAK
7      data_points: Dict[str, Any] # Supporting data
8      supports_hypothesis: bool   # True if supporting, False if
9      refuting
10
11 @dataclass

```

```

11  class EvidenceChain:
12      hypothesis_id: str
13      evidence_list: List[Evidence]
14      causal_pathway: List[str]      # A -> B -> C causal chain
15      overall_strength: EvidenceStrength
16      confidence_score: float      # 0.0 - 1.0
17      confidence_interval: Tuple[float, float] # (lower, upper)
18

```

Listing 5: Evidence-Based Hypothesis Generation

3.4 Hypothesis Structure

```

1  @dataclass
2  class DiagnosticHypothesis:
3      hypothesis_id: str
4      title: str                      # "High query volume causing
5      delays"
6      description: str
7      root_cause: str                  # Identified root cause
8      investigation_type: InvestigationType
9      entity_id: str                  # e.g., "SITE_001"
10     evidence_chain: EvidenceChain
11     recommendations: List[str]
12     alternative_hypotheses: List[str]
13     priority: str                  # "Critical", "High", "Medium", "
14     Low"
15
16     @property
17     def confidence(self) -> float:
18         return self.evidence_chain.confidence_score

```

Listing 6: Diagnostic Hypothesis

3.5 Confidence Calculation

```

1  def _recalculate_strength(self):
2      strength_map = {
3          EvidenceStrength.STRONG: 0.9,
4          EvidenceStrength.MODERATE: 0.6,
5          EvidenceStrength.WEAK: 0.3,
6      }
7
8      supporting = [e for e in self.evidence_list if e.
9      supports_hypothesis]
10     refuting = [e for e in self.evidence_list if not e.
11     supports_hypothesis]
12
13     support_score = sum(strength_map[e.strength] for e in
14     supporting)
15     refute_score = sum(strength_map[e.strength] for e in refuting)
16
17     # Net confidence with uncertainty bounds
18     self.confidence_score = max(0, (support_score - refute_score) /
19                               max(len(self.evidence_list), 1))

```

Listing 7: Evidence Strength to Confidence

The screenshot shows the Causal Hypothesis Engine interface with three separate hypothesis cards:

- Data Quality Anomaly**: Patient Study_22|Site_359|Subject_31763 in Study_22 shows critical DQI drop. Evidence Chain: Metric Outlier: DQI score of 33.3333206176758% is 3.5 SD from mean; Safety Signal: Related SAE pending review. Confidence: 0.95. Actions: Auto-Fix, Verify.
- Data Quality Anomaly**: Patient Study_21|Site_984|Subject_14164 in Study_21 shows critical DQI drop. Evidence Chain: Metric Outlier: DQI score of 66.66666412353516% is 3.5 SD from mean; Safety Signal: Related SAE pending review. Confidence: 0.95. Actions: Auto-Fix, Verify.
- Data Quality Anomaly**: Patient Study_21|Site_895|Subject_9459 in Study_21 shows critical DQI drop. Evidence Chain: (empty). Confidence: (empty). Actions: (empty).

Figure 4: Hypothesis Generation

4 Forecaster Agent (v1.0)

4.1 Purpose

The Forecaster Agent generates timeline predictions with uncertainty quantification, runs what-if simulations, and provides risk-adjusted forecasting for clinical trial milestones.

4.2 Forecast Types

Table 6: Forecaster Capabilities

Forecast	Output	Methodology
DB Lock	Date + CI	Issue resolution velocity
Clean Patient	Date + CI	DQI trajectory modeling
Issue Resolution	Days + CI	Historical pattern matching
Query Resolution	Date + CI	Query aging analysis
SDV Completion	Percentage + Date	Site performance benchmarks
Site Performance	DQI projection	Trend extrapolation

4.3 Uncertainty Quantification

```

1  @dataclass
2  class UncertaintyBand:
3      point_estimate: float          # Best estimate
4      lower_bound_95: float          # 95% CI lower
5      upper_bound_95: float          # 95% CI upper
6      lower_bound_80: float          # 80% CI lower
7      upper_bound_80: float          # 80% CI upper
8      lower_bound_50: float          # 50% CI lower
9      upper_bound_50: float          # 50% CI upper
10     standard_error: float
11     confidence_level: ConfidenceLevel # VERY_HIGH to VERY_LOW
12
13 class ConfidenceLevel(str, Enum):
14     VERY_HIGH = "very_high" # SE < 5%
15     HIGH = "high"          # SE 5-10%
16     MEDIUM = "medium"       # SE 10-20%
17     LOW = "low"            # SE 20-30%
18     VERY_LOW = "very_low"  # SE > 30%
19

```

Listing 8: Uncertainty Bands

4.4 What-If Scenarios

```

1  class ScenarioType(str, Enum):
2      ADD_RESOURCE = "add_resource"      # +1 CRA
3      REMOVE_RESOURCE = "remove_resource"
4      CLOSE_SITE = "close_site"
5      ADD_SITE = "add_site"
6      ACCELERATE_RESOLUTION = "accelerate_resolution"
7      DELAY = "delay"
8      PROCESS_CHANGE = "process_change"
9
10 @dataclass
11 class WhatIfScenario:
12     scenario_type: ScenarioType
13     description: str
14     parameters: Dict[str, Any]        # {"resource_count": 2}
15     baseline_outcome: Dict           # Without intervention
16     scenario_outcome: Dict          # With intervention
17     impact_delta: Dict              # Difference
18     confidence: float
19     recommendations: List[str]
20     cost_estimate: Optional[float]
21

```

Listing 9: Scenario Simulation

The screenshot displays the 'What-If Simulator' interface. At the top, there's a header with a lightning bolt icon and the text 'What-If Simulator' followed by the subtext 'Test operational decisions before execution'. Below this, there are two input fields: 'SCENARIO TYPE' set to 'Site Closure Impact' and 'TARGET ENTITY' set to 'US-001'. A large blue button labeled '▷ RUN SIMULATION' is centered below these fields.

Impact Analysis

TIMELINE DELAY +16 weeks (Regulatory Chain Lag)	SUBJECT TRANSFERS 289
TRANSFER SUCCESS_PROB 82%	ESTIMATED DROPOUTS 23-34 subjects
COST IMPACT \$449,000 (Transfer & Audit Cost)	

ALTERNATIVES ANALYZED

Close site	+6 weeks	No
Add coordinator	0 weeks	YES ✓
Increase monitoring	+2 weeks	Maybe

FINAL RECOMMENDATION
Site US-001 has 289 subjects. Adding a dedicated coordinator is 65% more cost-effective than closure.

APPROVE **REJECT**

Figure 5: What-If Scenario Simulation

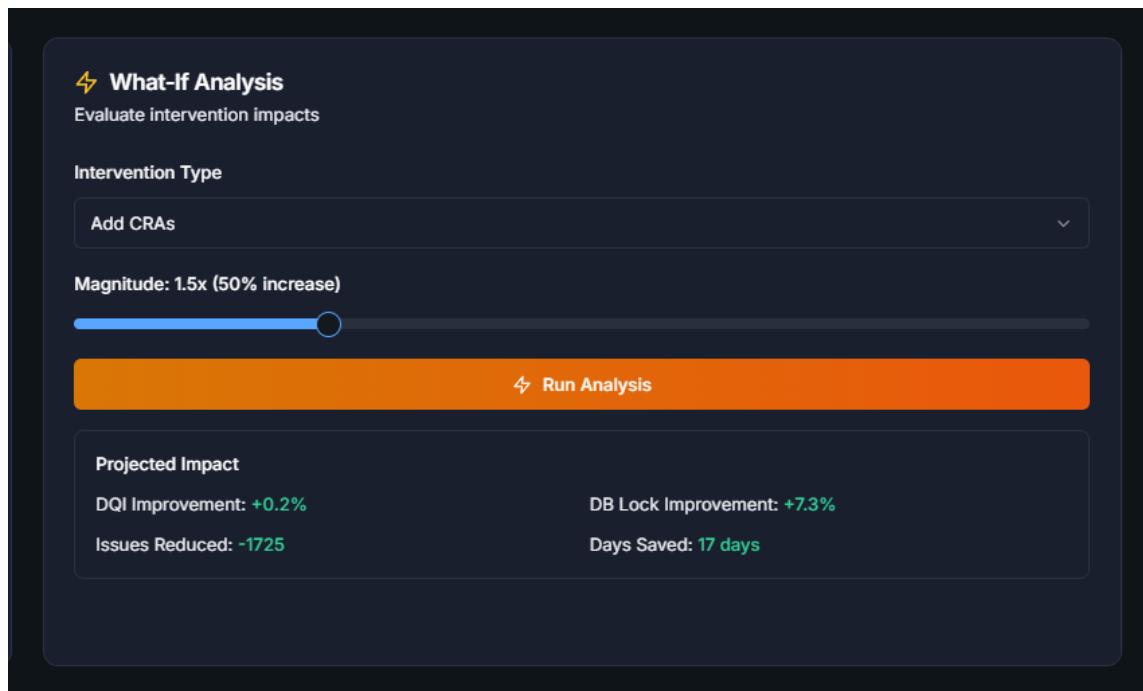


Figure 6: What-If Scenario Analysis

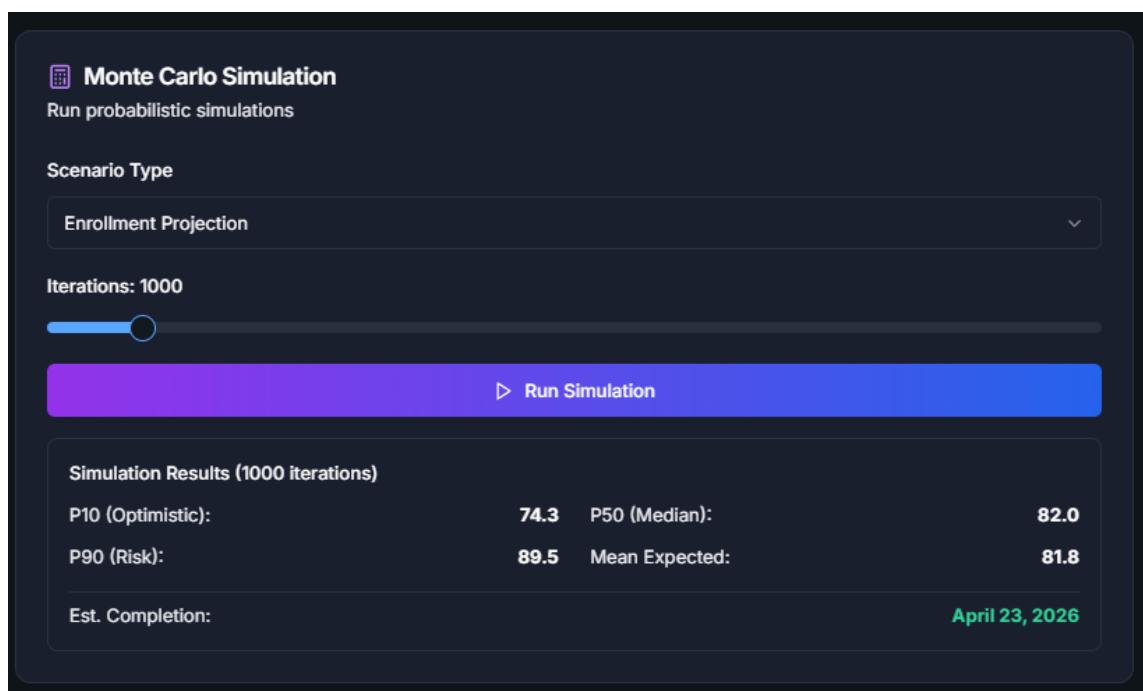


Figure 7: Monte-carlo Simulations

4.5 Timeline Milestones

```

1 @dataclass
2 class TimelineMilestone:
3     milestone_id: str
4     name: str
5     target_date: datetime
        # "Last Patient Last Visit"
        # Original target
    
```

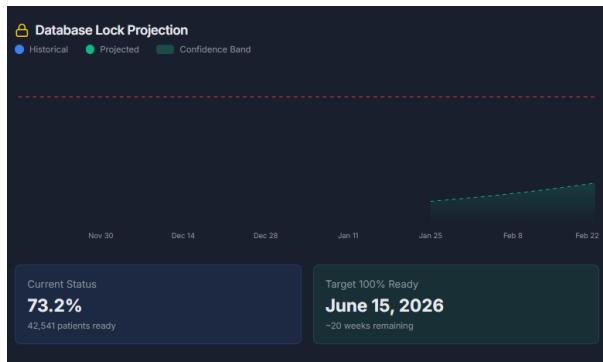


Figure 8: DB Lock Projection Timeline

```

6     predicted_date: datetime      # Current prediction
7     uncertainty: UncertaintyBand
8     probability_on_time: float    # P(actual <= target)
9     days_variance: int           # predicted - target
10    status: str                  # "on_track", "at_risk", "delayed"
11    blockers: List[str]          # ["High query volume at Site 05"]
12    accelerators: List[str]       # ["Dedicated CRA assigned"]
13

```

Listing 10: Milestone Tracking

5 Resolver Agent (v1.0)

5.1 Purpose

The Resolver Agent generates prioritized action plans, calculates cascade impact of resolutions, and provides role-based task assignments with effort estimates.

5.2 Action Plan Structure

```

1 @dataclass
2 class Action:
3     action_id: str
4     title: str
5     description: str
6     category: ActionCategory      # DATA_QUALITY, SAFETY,
MONITORING
7     priority: ActionPriority      # CRITICAL, HIGH, MEDIUM, LOW
8     responsible_role: ResponsibleRole
9     entity_id: str               # Target patient/site/study
10    estimated_effort_hours: float
11    dependencies: List[str]      # Other action IDs
12    due_date: Optional[datetime]
13    success_rate: float          # From resolution genome
14    requires_approval: bool
15
16 @dataclass
17 class ActionPlan:
18     plan_id: str
19     title: str
20     actions: List[Action]

```

```

21     total_effort_hours: float
22     expected_impact: Dict          # DQI improvement, etc.
23     cascade_impacts: List[CascadeImpact]
24

```

Listing 11: Action Plan Model

5.3 Role-Based Assignment

Table 7: Responsible Roles

Role	Scope	Action Types
CRA	Patient/Site	SDV, missing visits, signature gaps
Data Manager	Query/Data	Query resolution, data fixes
Safety Data Manager	SAE	SAE-DM pending, safety queries
Safety Physician	SAE-Safety	SAE-Safety pending
Medical Coder	Coding	MedDRA/WHODrug uncoded
Site Coordinator	Site	Training, process changes
Study Lead	Study	Escalations, resource decisions

5.4 Cascade Impact Analysis

```

1  @dataclass
2  class CascadeImpact:
3      source_issue: str           # Issue being resolved
4      source_entity: str         # Patient/site
5      direct_impact: Dict        # Immediate effect
6      cascade_effects: List[Dict] # Downstream effects
7      total_issues_resolved: int # Including cascade
8      total_patients_unblocked: int
9      dqi_improvement: float    # Expected DQI gain
10     db_lock_acceleration_days: float
11     effort_hours: float
12     roi_score: float          # Impact / effort
13
14     # ROI-based prioritization
15     def prioritize_actions(actions: List[Action]) -> List[Action]:
16         for action in actions:
17             cascade = calculate_cascade_impact(action)
18             action.roi_score = cascade.total_patients_unblocked /
action.effort_hours
19         return sorted(actions, key=lambda a: a.roi_score, reverse=True)
20

```

Listing 12: Cascade Impact Calculation

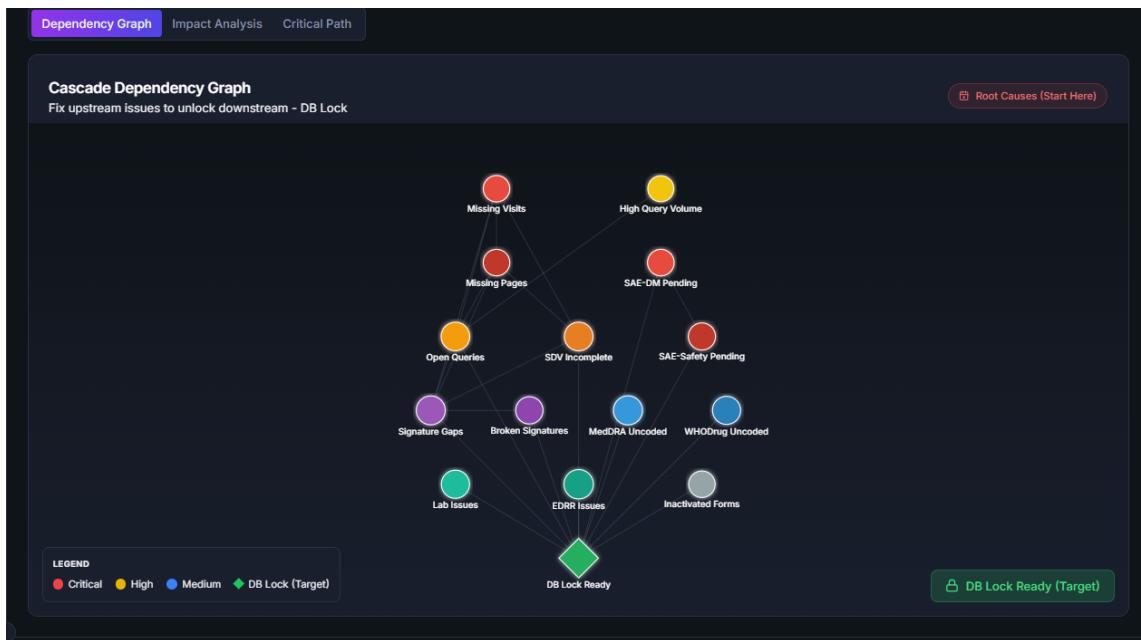


Figure 9: Cascade Impact Analysis Graph

6 Executor Agent (v1.0)

6.1 Purpose

The Executor Agent validates actions, manages approval workflows, executes with audit trails, and provides rollback capabilities for failed executions.

6.2 Execution Status Flow

```

1  class ExecutionStatus(str, Enum):
2      PENDING = "pending"
3      VALIDATING = "validating"
4      VALIDATION_FAILED = "validation_failed"
5      AWAITING_APPROVAL = "awaiting_approval"
6      APPROVED = "approved"
7      REJECTED = "rejected"
8      IN_PROGRESS = "in_progress"
9      COMPLETED = "completed"
10     FAILED = "failed"
11     ROLLED_BACK = "rolled_back"
12     EXPIRED = "expired"
13

```

Listing 13: Execution Status Lifecycle

6.3 Validation Framework

```

1  @dataclass
2  class ValidationCheck:
3      check_id: str
4      name: str          # "data_integrity"
5      description: str

```

```

1      result: ValidationResult      # PASSED, WARNING, FAILED
2      message: str
3      details: Dict[str, Any]
4
5  @dataclass
6  class ValidationReport:
7      action_id: str
8      checks: List[ValidationCheck]
9      overall_result: ValidationResult
10     risk_level: RiskLevel          # LOW, MEDIUM, HIGH, CRITICAL
11     can_proceed: bool              # True only if overall_result != FAILED
12     approval_required: ApprovalLevel
13
14
15
16
17
18

```

Listing 14: Pre-Execution Validation

6.4 Approval Workflow

Table 8: Approval Levels by Risk

Risk Level	Approval	Approvers	Timeout
Low	None	0	—
Medium	Single	1	24h
High	Dual	2	12h
Critical	Study Lead	1+	4h

6.5 Audit Trail

```

1  @dataclass
2  class AuditEntry:
3      entry_id: str
4      action_id: str
5      event_type: str          # "validation", "approval", "execution"
6      event_description: str
7      actor: str                # User or system ID
8      actor_role: str
9      old_state: Dict[str, Any]
10     new_state: Dict[str, Any]
11     timestamp: datetime
12     checksum: str            # Tamper detection
13
14     def __post_init__(self):
15         # Calculate checksum for audit integrity
16         self.checksum = self._calculate_checksum()
17
18     def _calculate_checksum(self) -> str:
19         content = f"{self.action_id}|{self.event_type}|{self.actor}"
20         return hashlib.sha256(content.encode()).hexdigest()[:16]
21

```

Listing 15: 21 CFR Part 11 Compliant Audit

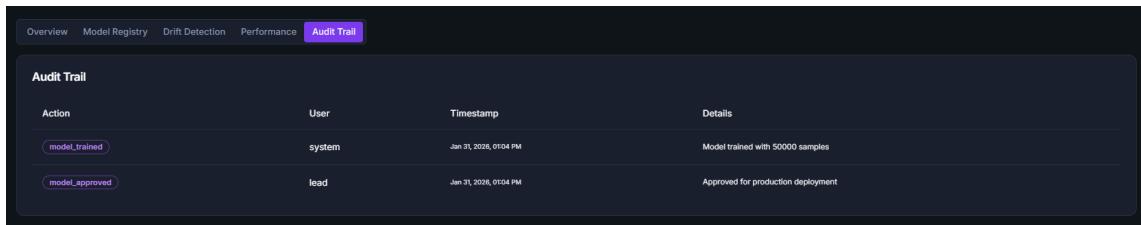


Figure 10: Audit Trail

7 Communicator Agent (v1.0)

7.1 Purpose

The Communicator Agent drafts context-aware messages, selects appropriate channels based on recipient preferences, and handles notification batching.

7.2 Message Types

Table 9: Message Types and Channels

Type	Priority	Channel	Approval
Alert	Urgent	SMS + Email	No
Escalation	High	Email + In-App	Yes
Report	Normal	Email	No
Notification	Normal	In-App	No
Reminder	Low	In-App	No
Summary	Normal	Email (batch)	No

7.3 Recipient Profiles

```

1 @dataclass
2 class RecipientProfile:
3     recipient_id: str
4     name: str
5     email: str
6     role: RecipientRole
7     preferred_channel: Channel      # EMAIL, SMS, IN_APP, SLACK
8     preferred_language: str         # "en", "es", "fr"
9     timezone: str                  # For scheduling
10    quiet_hours: Tuple[int, int]   # (22, 7) = 10PM to 7AM
11    batch_notifications: bool      # Combine into digest
12    max_daily_messages: int       # Rate limiting
13    sites: List[str]              # Subscribed sites
14    studies: List[str]            # Subscribed studies
15

```

Listing 16: Recipient Preference Management

7.4 Template System

```

1  @dataclass
2  class MessageTemplate:
3      template_id: str
4      name: str
5      message_type: MessageType
6      subject_template: str
7      body_template: str
8      variables: List[str]
9      requires_approval: bool
10
11     def render(self, context: Dict[str, Any]) -> Tuple[str, str]:
12         """Render template with context variables."""
13         subject = self.subject_template.format(**context)
14         body = self.body_template.format(**context)
15         return subject, body
16

```

Listing 17: Message Template Rendering

7.5 Notification Batching

```

1  @dataclass
2  class NotificationBatch:
3      batch_id: str
4      recipient: RecipientProfile
5      messages: List[Message]
6      batch_type: str
7      scheduled_time: datetime
8      status: DeliveryStatus
9
10     def add_message(self, message: Message):
11         self.messages.append(message)
12
13     @property
14     def total_messages(self) -> int:
15         return len(self.messages)
16

```

Listing 18: Daily Digest Batching

8 Agent Tools & Memory

8.1 Tool Registry

The ToolRegistry provides 50+ tools organized by category:

Table 10: Tool Categories

Category	Purpose	Tools
data	Patient/site/study retrieval	10+
search	Vector search, RAG knowledge	5
analytics	DQI, cascade, benchmarks	12+
action	Create/update actions	8
ml	Model predictions	6
cross-study	Pattern transfer	6+

8.2 Key Data Tools

```

1 # Patient retrieval
2 get_patient(patient_key: str) -> ToolResult
3 get_high_priority_patients(limit: int = 20) -> ToolResult
4
5 # Site analytics
6 get_site_summary(site_id: str) -> ToolResult
7 get_lowest_dqi_sites(limit: int = 10) -> ToolResult
8 get_site_patient_details(site_id: str) -> ToolResult
9
10 # Study analytics
11 get_study_summary(study_id: str) -> ToolResult
12 get_overall_summary() -> ToolResult
13

```

Listing 19: Core Data Tools

8.3 Analytics Tools

```

1 # Cascade analysis
2 get_cascade_impact(patient_key: str) -> ToolResult
3
4 # DQI analysis
5 get_dqi_distribution(entity_type: str, entity_id: str) ->
ToolResult
6
7 # Benchmarking
8 get_site_benchmarks(site_id: str) -> ToolResult
9
10 # Pattern detection
11 find_similar_patterns(study_id: str, pattern_type: str) ->
ToolResult
12     get_pattern_transfer_recommendations(source: str, target: str) ->
ToolResult
13

```

Listing 20: Analytics Tools

8.4 Memory System

The memory system enables agents to learn across sessions:

Table 11: Memory Modules

Module	File	Purpose
Conversation	<code>conversation_memory.py</code>	Session context retention
Decision History	<code>decision_history.py</code>	Past routing decisions
Learning Engine	<code>learning_engine.py</code>	Pattern extraction
Persistent Memory	<code>persistent_memory.py</code>	PostgreSQL long-term storage

8.5 Persistent Memory

```

1  @dataclass
2  class Memory:
3      memory_id: str
4      agent_id: str
5      context: str          # When this learning applies
6      learning: str         # The actual insight
7      confidence: float    # 0-1
8      access_count: int    # Usage frequency
9      relevance_score: float # Semantic similarity
10
11 class PersistentAgentMemory:
12     def remember(self, context: str, learning: str, confidence: float):
13         """Store a learning in PostgreSQL."""
14         memory = Memory(
15             memory_id=str(uuid.uuid4()),
16             agent_id=self.agent_id,
17             context=context,
18             learning=learning,
19             confidence=confidence
20         )
21         self._persist_memory(memory)
22
23     def recall(self, context: str, limit: int = 5) -> List[Memory]:
24         """Retrieve relevant memories using semantic search."""
25         query_embedding = self.embedding_model.encode(context)
26         return self._find_similar(query_embedding, limit)
27
28     def consolidate(self, min_confidence: float = 0.3):
29         """Remove low-confidence old memories."""
30         # Keep memories with high access or high confidence
31

```

Listing 21: Cross-Session Learning

9 Data Loaders

Each agent has a specialized data loader for efficient data access:

```

1  class DiagnosticDataLoader:
2      def __init__(self, base_path: str = "data/processed"):
3          self.base_path = Path(base_path)
4          self._cache: Dict[str, pd.DataFrame] = {}
5
6      # Cached property loaders

```

```

7     @property
8     def patient_issues(self) -> pd.DataFrame:
9         return self._load_parquet("patient_issues.parquet")
10
11    @property
12    def patient_cascade(self) -> pd.DataFrame:
13        return self._load_parquet("cascade/patient_cascade.parquet")
14
15    # Entity-specific retrievers
16    def get_patient_data(self, patient_key: str) -> Dict:
17        """Get all data for a specific patient."""
18
19    def get_site_data(self, site_id: str) -> Dict:
20        """Get aggregated data for a site."""
21
22    def get_study_data(self, study_id: str) -> Dict:
23        """Get study-level aggregations."""
24

```

Listing 22: Agent Data Loader Pattern

10 Report Generation

10.1 Overview

The Report Generation module (`src/generation/`) provides automated PDF, Word, and PowerPoint report generation with professional styling and multi-channel delivery.

Table 12: Report Generation Components

Component	File	Purpose
Report Generators	<code>report_generators.py</code>	8 report types with templates
Export Engine	<code>export_engine.py</code>	PDF/DOCX/PPTX export
Template Engine	<code>template_engine.py</code>	Jinja2-based templating
Scheduler	<code>report_scheduler.py</code>	Cron-based automation
Auto Summarizer	<code>auto_summarizer.py</code>	LLM-powered summaries
NL Interface	<code>nl_interface.py</code>	Natural language queries

10.2 Report Types

Table 13: Available Report Types

Report Type	Description	Formats
Portfolio Summary	Overall portfolio status	PDF, DOCX
Study Summary	Study-level KPIs	PDF, DOCX
Site Performance	Site benchmarks	PDF, DOCX
Patient Status	Individual patient details	PDF
Risk Assessment	Risk tier breakdown	PDF, PPTX
DQI Analysis	Data Quality metrics	PDF, DOCX
Issue Report	Open issues summary	PDF, DOCX
Executive Summary	Leadership briefing	PDF, PPTX

10.3 Export Engine

```
1 class OutputFormat:
2     HTML = "html"
3     PDF = "pdf"
4     DOCX = "docx"
5     PPTX = "pptx"
6     CSV = "csv"
7     JSON = "json"
8
9 @dataclass
10 class StyleConfig:
11     primary_color: str = "#1a365d"
12     secondary_color: str = "#2c5282"
13     accent_color: str = "#38a169"
14     warning_color: str = "#dd6b20"
15     danger_color: str = "#e53e3e"
16     font_family: str = "Arial, Helvetica, sans-serif"
17     heading1_size: int = 18
18     body_size: int = 11
19     page_margin_inches: float = 1.0
20     company_name: str = "TrialPlus Nexus"
21
22 class PDFExporter:
23     """PDF export with multiple backend fallbacks."""
24     BACKENDS = ["weasyprint", "pdfkit", "xhtml2pdf", "reportlab"]
25
26     def export(self, html_content: str, output_path: str) ->
27         ExportResult:
28         # Try backends in priority order
29         for backend in self.available_backends:
30             try:
31                 return self._export_with_backend(backend,
32                     html_content, output_path)
33             except Exception as e:
34                 continue
35         return ExportResult(success=False, error="All backends
36 failed")
37
38
39
```

Listing 23: Export Format Configuration

10.4 Data Loading

```
1 class DataLoader:
2     def __init__(self):
3         self.data_dir = Path("data/processed")
4         self._cache: Dict[str, pd.DataFrame] = {}
5
6     def get_portfolio_summary(self) -> Dict:
7         """Get portfolio-level summary statistics."""
8         return {
9             'total_patients': self._count_patients(),
10            'total_sites': self._count_sites(),
11            'overall_dqi': self._calculate_portfolio_dqi(),
12            'high_risk_count': self._count_by_risk('high'),
13            'clean_patient_rate': self._clean_patient_rate(),
```

```

14         'db_lock_readiness': self._db_lock_readiness()
15     }
16
17     def get_site_summary(self, site_id: str) -> Dict:
18         """Get site-level summary with benchmarks."""
19
20     def get_patient_data(self, site_id: str = None) -> pd.DataFrame
21     :
22         """Get unified patient record with optional filtering."""

```

Listing 24: Report Data Loader

10.5 Report Scheduler

```

1  class ReportFrequency(str, Enum):
2      DAILY = "daily"
3      WEEKLY = "weekly"
4      BIWEEKLY = "biweekly"
5      MONTHLY = "monthly"
6      ON_DEMAND = "on_demand"
7
8  class ReportScheduler:
9      def schedule_report(
10          self,
11          report_type: str,
12          frequency: ReportFrequency,
13          recipients: List[str],
14          delivery_channels: List[str] # ["email", "slack"]
15      ):
16          """Schedule a recurring report generation."""
17
18      def trigger_immediate(self, report_type: str, params: Dict):
19          """Generate and deliver a report immediately."""
20

```

Listing 25: Scheduled Report Generation

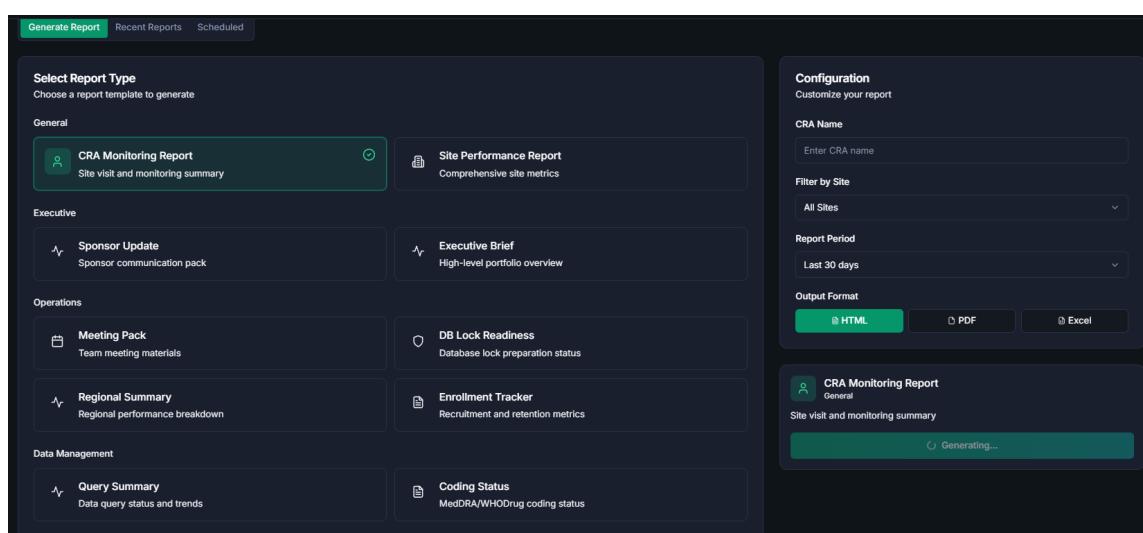


Figure 11: Report Generation Dashboard

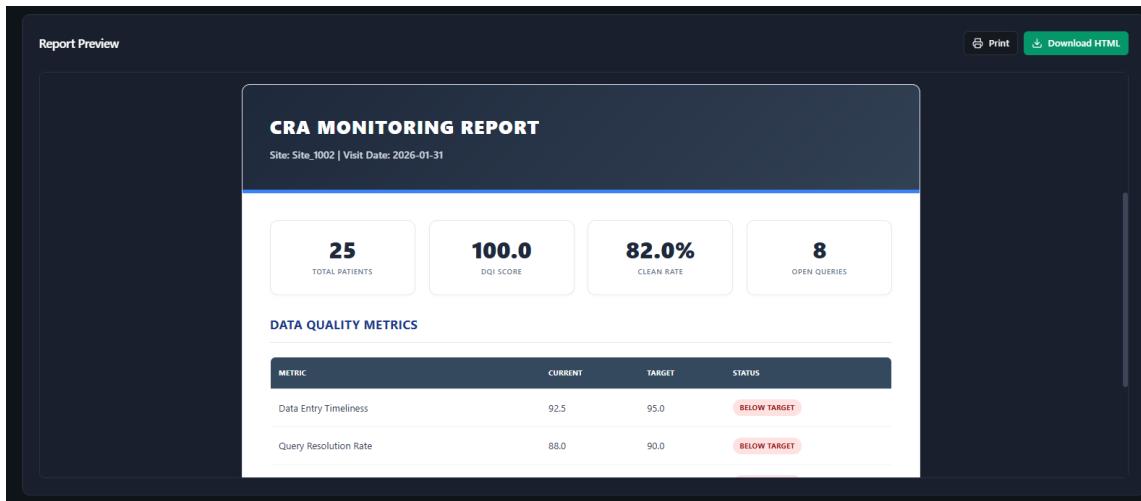


Figure 12: Sample Report Preview

11 Knowledge & Intelligence

11.1 Overview

The Knowledge module (`src/knowledge/`) provides RAG-based retrieval, graph intelligence, and causal analysis capabilities used by agent tools.

Table 14: Knowledge Components

Component	File	Used By
RAG Knowledge Base	<code>rag_knowledge_base.py</code>	NL Interface, Document Engine
Neo4j Graph	<code>neo4j_graph.py</code>	Agent Tools, Graph Routes
Causal Hypothesis	<code>causal_hypothesis_engine.py</code>	Agent Tools, Intelligence API
Hybrid Search	<code>hybrid_search.py</code>	RAG retrieval

11.2 RAG Knowledge Base

```

1  class DocumentChunker:
2      def __init__(self, chunk_size: int = 500, chunk_overlap: int =
3          100):
4          self.separators = ['\n\n', '\n', '.', '']
5
6          def chunk_text(self, text: str, metadata: Dict = None) -> List[Dict]:
7              """Chunk text into smaller pieces with metadata."""
8
9          class RAGKnowledgeBase:
10              def __init__(self, embedding_model: str = "all-MiniLM-L6-v2"):
11                  self.model = SentenceTransformer(embedding_model)
12                  self.chunker = DocumentChunker()
13
14              def load_all_documents(self) -> Dict[str, int]:
15                  """Load ICH-GCP, Protocol, and SOP documents."""
16                  return {
17                      'ich_gcp': self._load_ich_gcp(),
18                      'protocol': self._load_protocol_knowledge(),
19                  }

```

```

18         'sop': self._load_sops()
19     }
20
21     def search(self, query: str, top_k: int = 5) -> List[Dict]:
22         """Semantic search over knowledge base."""
23         query_embedding = self.model.encode([query])[0]
24         scores = np.dot(self.embeddings, query_embedding)
25         return self._get_top_results(scores, top_k)
26

```

Listing 26: RAG Document Chunking

11.3 Neo4j Knowledge Graph

```

1  class Neo4jGraphService:
2      """Entity graph for Cascade Intelligence."""
3
4      # Node Types: Study -> Site -> Patient -> Issue
5      # Relationships: CONTAINS, HAS_PATIENT, HAS_ISSUE, BLOCKS
6
7      def build_entity_graph(self, upr_df: pd.DataFrame):
8          """Build graph from Unified Patient Record."""
9          # Creates Study, Site, Patient, Issue nodes
10         # Creates all relationships
11
12         def get_cascade_path(self, issue_id: str, max_depth: int = 3):
13             """Compute cascade impact path from an issue."""
14             query = """
15                 MATCH path = (issue:Issue {issue_id: $issue_id}) -[*1..%
16                 d]-(related)
17                 RETURN path, [rel in relationships(path) | type(rel)]
18             as rel_types
19             """
20
21             # Returns affected patients, downstream dependencies
22

```

Listing 27: Entity Graph Service

11.4 Causal Hypothesis Engine

```

1  CAUSAL_PATTERNS = {
2      "PI_ABSENCE": {
3          "description": "Principal Investigator absence pattern",
4          "indicators": ["high missing_pi_signature", "consent_issues"]
5      },
6      "STAFF_TURNOVER": {
7          "description": "Site staff turnover pattern",
8          "indicators": ["sudden DQI drop", "training_incomplete"]
9      },
10     "PROTOCOL_COMPLEXITY": {
11         "description": "Protocol complexity causing errors",
12         "indicators": ["high query rates", "visit deviations"]
13     }
14 }

```

```
16     class CausalHypothesisEngine:  
17         def analyze_population(self, sample_size: int = 20) -> List[  
18             Hypothesis]:  
19             """Generate hypotheses with evidence chains."""  
20             # Analyzes site/patient data against known patterns  
21             # Returns prioritized hypotheses with confidence scores
```

Listing 28: Hypothesis Generation

TrialPlus Nexus 10X

Technical Documentation

Module 5: User Interface & Dashboards

Role-Based Views, Collaboration & Visual Analytics

Version: 1.0.0

Date: January 2026

Views: 8 Role-Based Dashboards

Framework: React + TypeScript + Vite

Clinical Trial Risk Intelligence Engine

Contents

1 UI Architecture Overview	3
1.1 Technology Stack	3
1.2 Application Structure	3
1.3 Navigation & Routing	3
2 Study Lead Dashboard	3
2.1 Purpose	3
2.2 Key Features	3
2.3 Data Sources	4
3 CRA (Clinical Research Associate) View	5
3.1 Purpose	5
3.2 Key Features	5
3.3 Workflow Support	6
4 Data Manager Hub	7
4.1 Purpose	7
4.2 Key Features	7
4.3 Key Metrics Displayed	8
5 Safety View	9
5.1 Purpose	9
5.2 Key Features	9
5.3 Regulatory Timelines	9
6 Medical Coder View	10
6.1 Purpose	10
6.2 Key Features	10
6.3 Coding Workflow	11
7 Site Portal	12
7.1 Purpose	12
7.2 Key Features	12
7.3 Site Performance Indicators	13
8 Executive Overview	14
8.1 Purpose	14
8.2 Key Features	14
9 Collaboration Hub	15
9.1 Purpose	15
9.2 Key Features	15
9.3 Communication Channels	16
10 Cascade Explorer	16
10.1 Purpose	16
10.2 Key Features	16
10.3 Graph Visualization	17

11 AI Assistant	17
11.1 Purpose	17
11.2 Key Features	17
11.3 Example Queries	18
12 Reports & Export	18
12.1 Purpose	18
12.2 Available Reports	19
12.3 Report Generation Workflow	19
13 ML Governance Dashboard	19
13.1 Purpose	19
13.2 Key Features	19
13.3 Model Lifecycle States	20
14 Visualization Dashboard	21
14.1 Purpose	21
14.2 Visualization Types	21

1 UI Architecture Overview

1.1 Technology Stack

The frontend is built using modern React with TypeScript, providing type-safe development and optimal performance.

Table 1: Frontend Technology Stack

Technology	Purpose	Version
React 18	Component framework	18.x
TypeScript	Type safety	5.x
Vite	Build tool & dev server	Latest
Zustand	State management	Latest
Recharts	Data visualization	Latest
Axios	API communication	Latest

1.2 Application Structure

The application follows a feature-based architecture where each role-based view is a self-contained module:

Table 2: Frontend Directory Structure

Directory	Contents
src/features/	Role-based view components (8 views)
src/components/	Shared UI components (18 components)
src/services/	API service layer
src/stores/	Zustand state stores
src/types/	TypeScript type definitions

1.3 Navigation & Routing

Users navigate between role-based views through a sidebar that adapts to their assigned role. Each view is optimized for specific clinical operations workflows.

2 Study Lead Dashboard

2.1 Purpose

The Study Lead Dashboard provides a comprehensive portfolio overview for clinical operations leadership. It is the highest-level view showing cross-study performance and strategic insights.

2.2 Key Features

- Portfolio Summary Cards:** Total patients, sites, studies, and overall DQI score displayed prominently at the top
- DB Lock Readiness Gauge:** Visual gauge showing percentage of patients ready for database lock
- Clean Patient Progress:** Donut chart showing clean vs. pending patients

- **Risk Distribution:** Breakdown of patients by risk tier (Critical, High, Medium, Low)
- **Site Performance Table:** Sortable table ranking sites by performance metrics
- **AI Assistant:** Integrated chat interface for natural language queries about study data

2.3 Data Sources

The Study Lead view aggregates data from multiple API endpoints:

- `/dashboards/study_lead` – Primary dashboard metrics
- `/analytics/portfolio` – Portfolio-level statistics
- `/intelligence/hypotheses` – AI-generated insights
- `/simulation/db-lock-projection` – Timeline projections

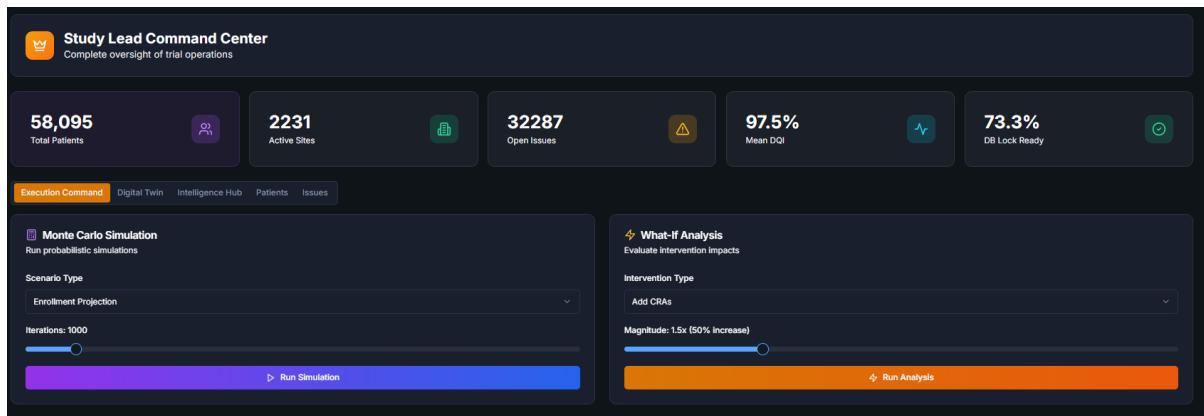


Figure 1: Study Lead View 1: Overview

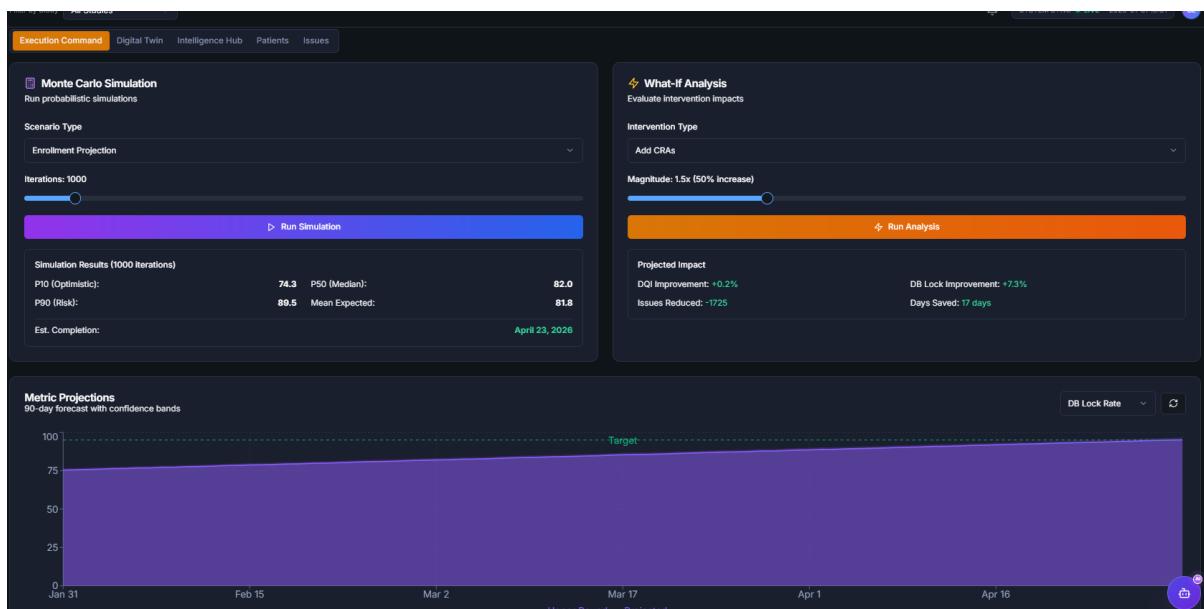


Figure 2: Study Lead View 2: Monte-Carlo Simulation & What-If Analysis

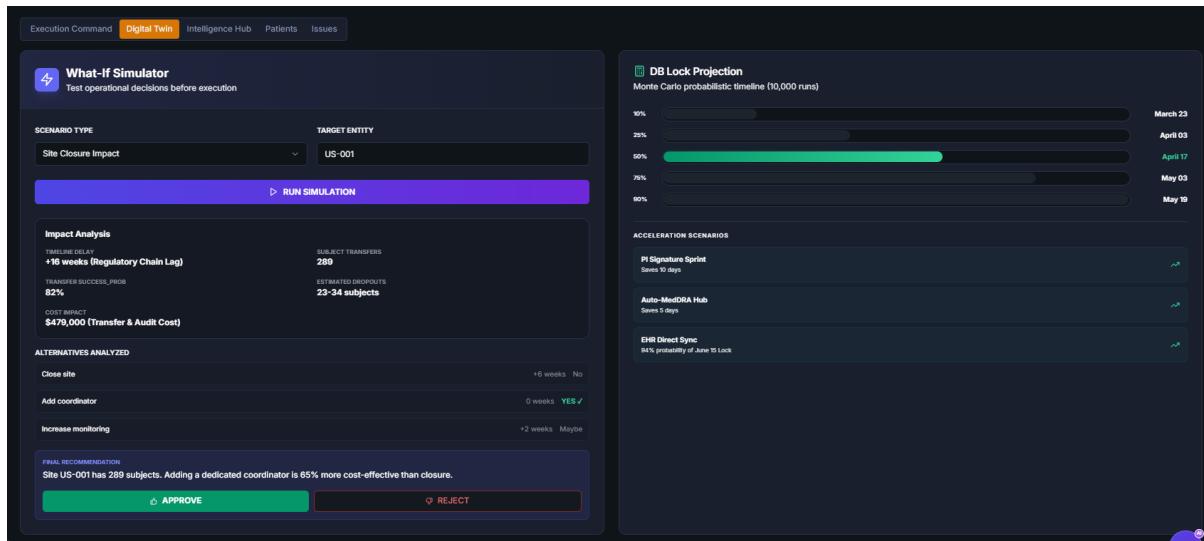


Figure 3: Study Lead View 3: What-If Simulator

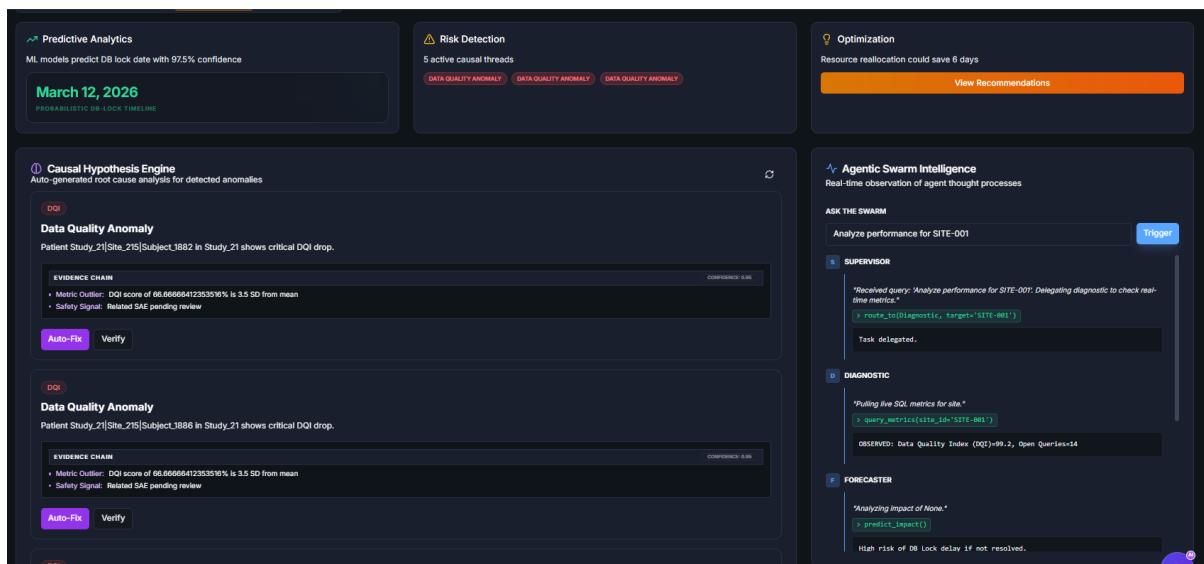


Figure 4: Study Lead View 4: Causal Hypothesis Engine

3 CRA (Clinical Research Associate) View

3.1 Purpose

The CRA View is designed for field monitors who need to track patient issues, plan site visits, and manage their monitoring workload efficiently.

3.2 Key Features

- Smart Work Queue:** AI-prioritized list of actions based on impact and urgency
- Site Summary Cards:** Quick view of assigned sites with key metrics

- **Patient Issues Table:** Filterable list of open issues requiring attention
- **Issue Resolution Workflow:** One-click resolution with audit trail
- **SDV Progress Tracker:** Source Data Verification completion status
- **Visit Planning Calendar:** Recommended visit schedule based on site priorities

3.3 Workflow Support

The CRA view supports the complete monitoring workflow:

1. Review prioritized work queue at start of day
2. Select patient or site to investigate
3. View detailed issue breakdown with AI recommendations
4. Resolve issues with documented reason for change
5. Track resolution impact on site DQI

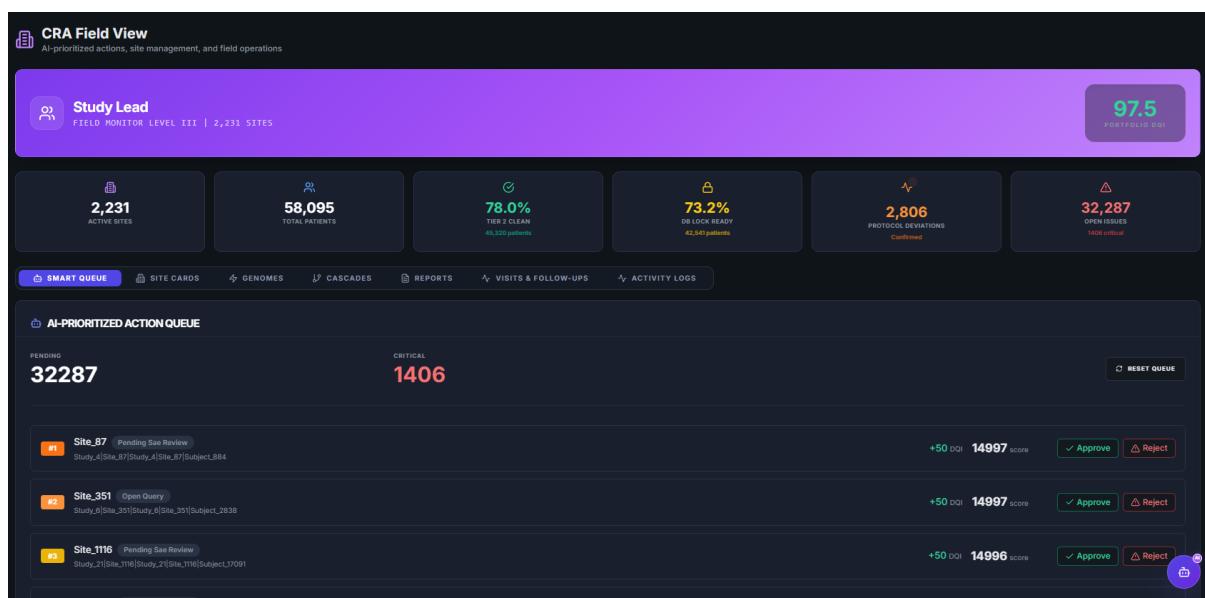


Figure 5: CRA View 1: Overview

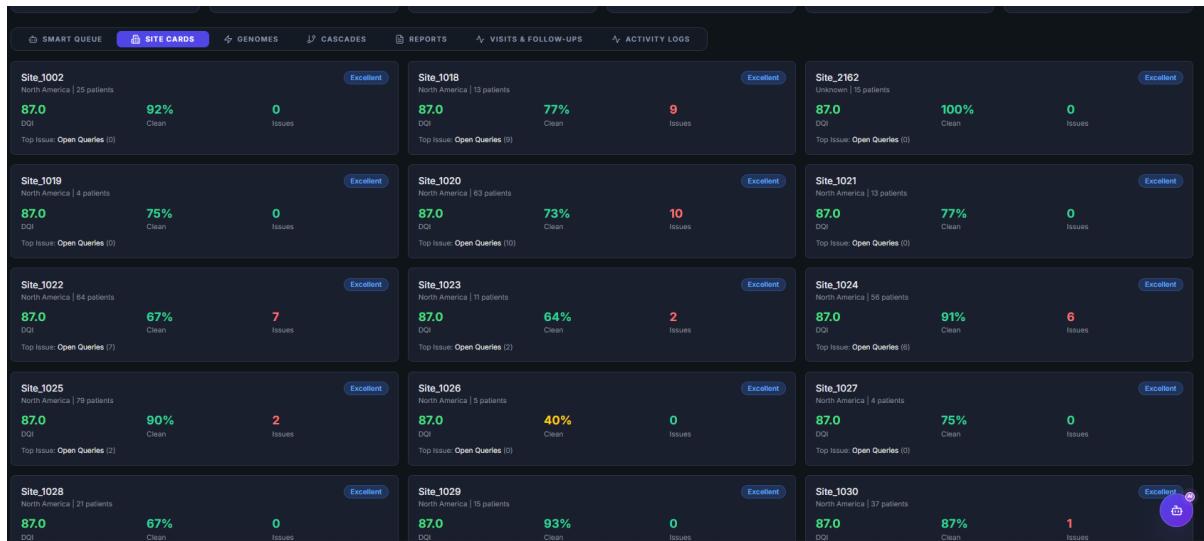


Figure 6: CRA View 2: Site Cards

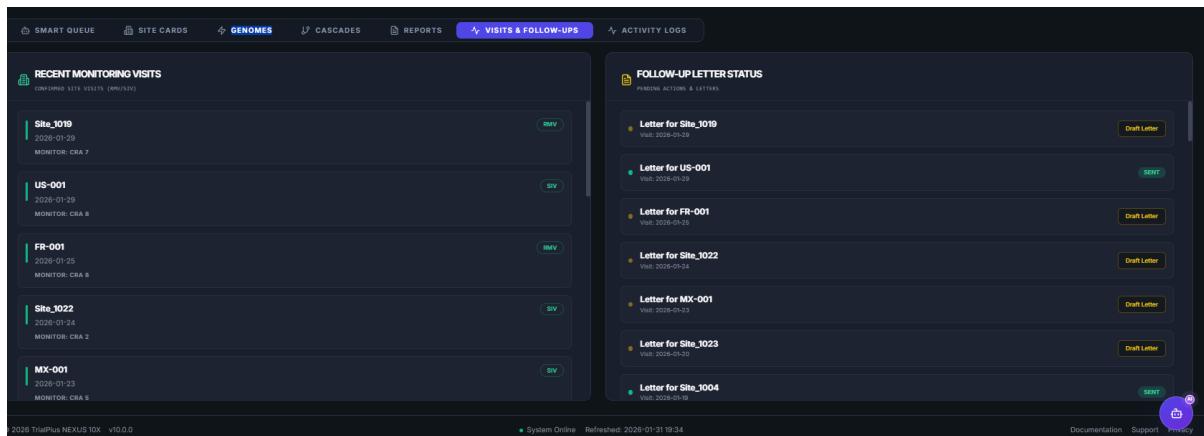


Figure 7: CRA View 3: Visits & Follow-ups

4 Data Manager Hub

4.1 Purpose

The Data Manager Hub (DMHub) centralizes all data quality management activities, including query management, edit checks, and data reconciliation.

4.2 Key Features

- **Query Aging Dashboard:** Visual breakdown of query age distribution
- **Issue Summary Metrics:** Total open, resolved, and pending issues
- **DQI Trend Chart:** Historical DQI score progression
- **Bottleneck Analysis:** Identification of systemic data issues
- **LPLV Tracker:** Last Patient Last Visit status monitoring

- **Bulk Resolution Tools:** Batch processing for similar issues

4.3 Key Metrics Displayed

Table 3: Data Manager Dashboard Metrics

Metric	Type	Purpose
Open Queries	Count	Workload indicator
Query Resolution Rate	Percentage	Efficiency metric
Average Query Age	Days	SLA monitoring
DQI Score	0-100	Data quality health
DB Lock Ready	Percentage	Milestone readiness

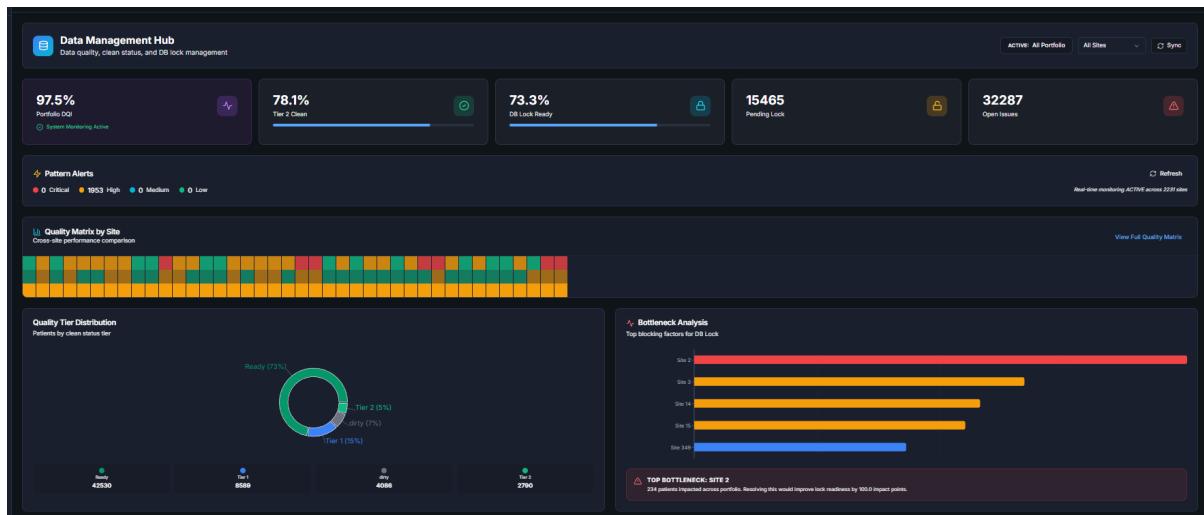


Figure 8: DM Hub View 1: Overview

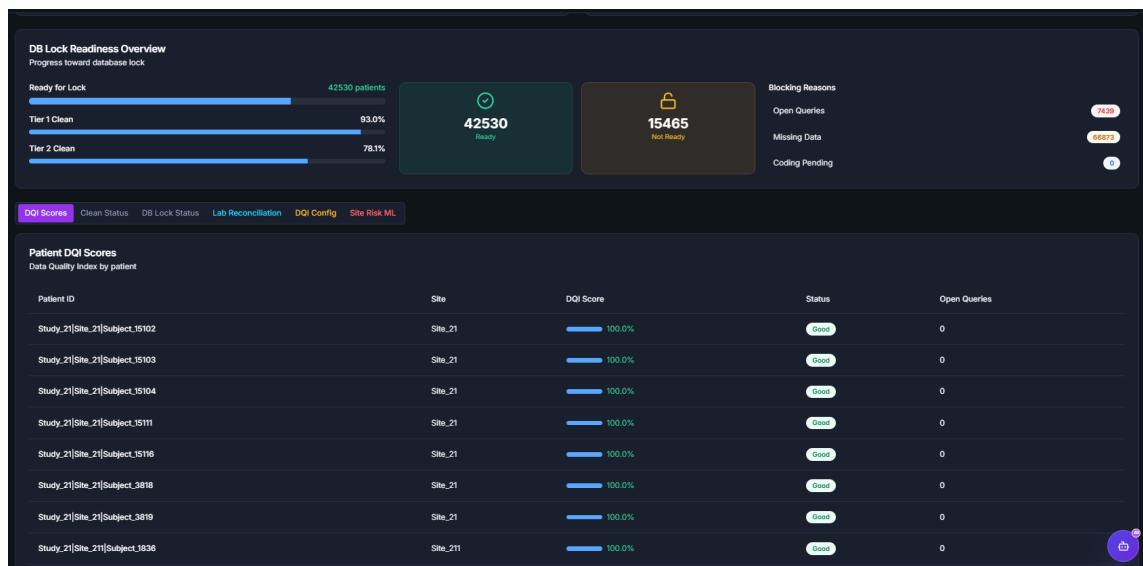


Figure 9: DM Hub View 2: DB Lock Readiness & DQI Scores

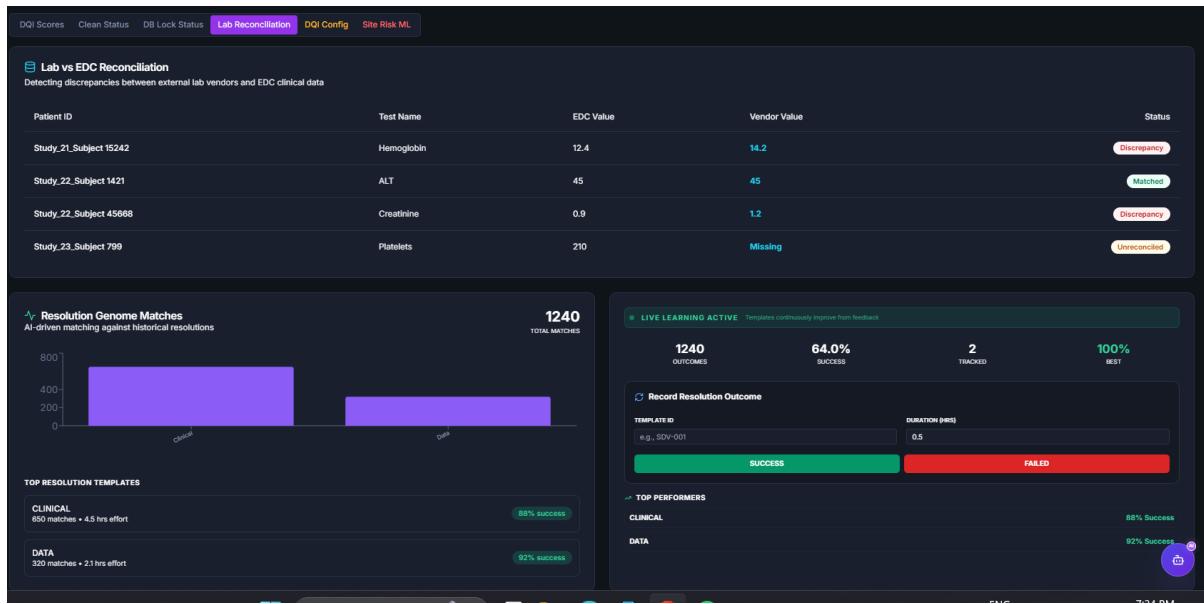


Figure 10: DM Hub View 3: Lab Reconciliation & Resolution Genome

5 Safety View

5.1 Purpose

The Safety View is designed for pharmacovigilance teams to monitor Serious Adverse Events (SAEs), track regulatory timelines, and detect safety signals.

5.2 Key Features

- SAE Case List:** Tabular view of all SAE cases with status and timeline
- SLA Countdown:** Visual indicators for regulatory reporting deadlines
- Safety Signal Detection:** Flagged patterns requiring medical review
- Narrative Generation:** AI-assisted safety narrative drafting
- Timeline Visualization:** SAE progression over time
- Pending Actions:** DM-pending and Safety-pending case tracking

5.3 Regulatory Timelines

The Safety View enforces critical regulatory deadlines:

Table 4: Safety Reporting SLAs

Event Type	Initial Report	Follow-up
Fatal/Life-threatening	7 calendar days	15 calendar days
Other SAEs	15 calendar days	As needed
SUSARs	7-15 days	Per regulation

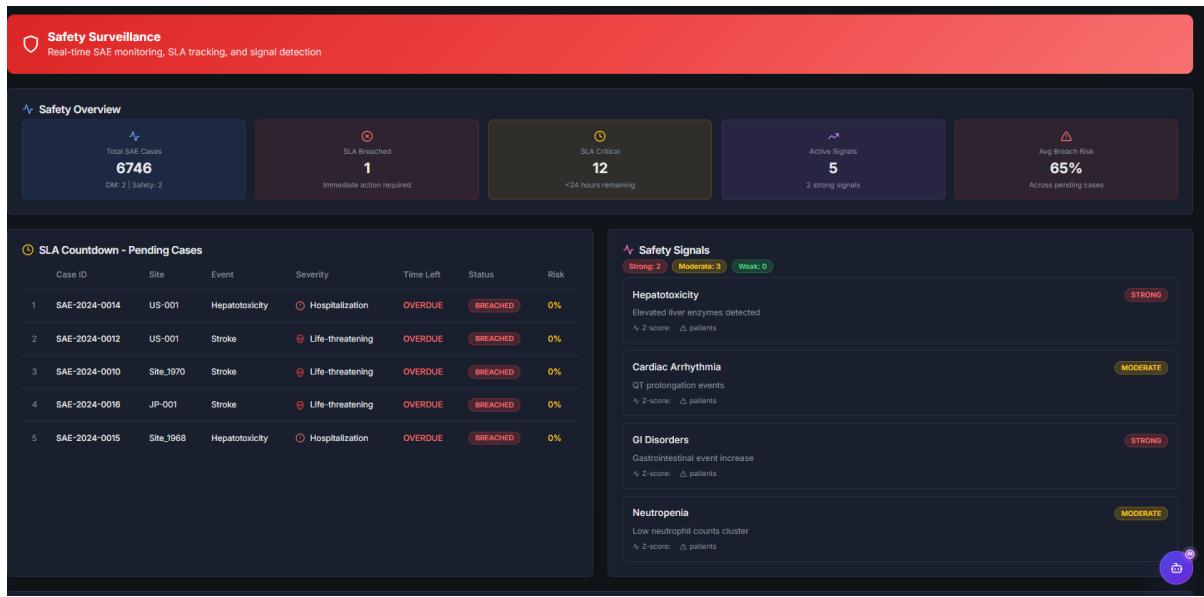


Figure 11: Safety Surveillance View 1: Overview

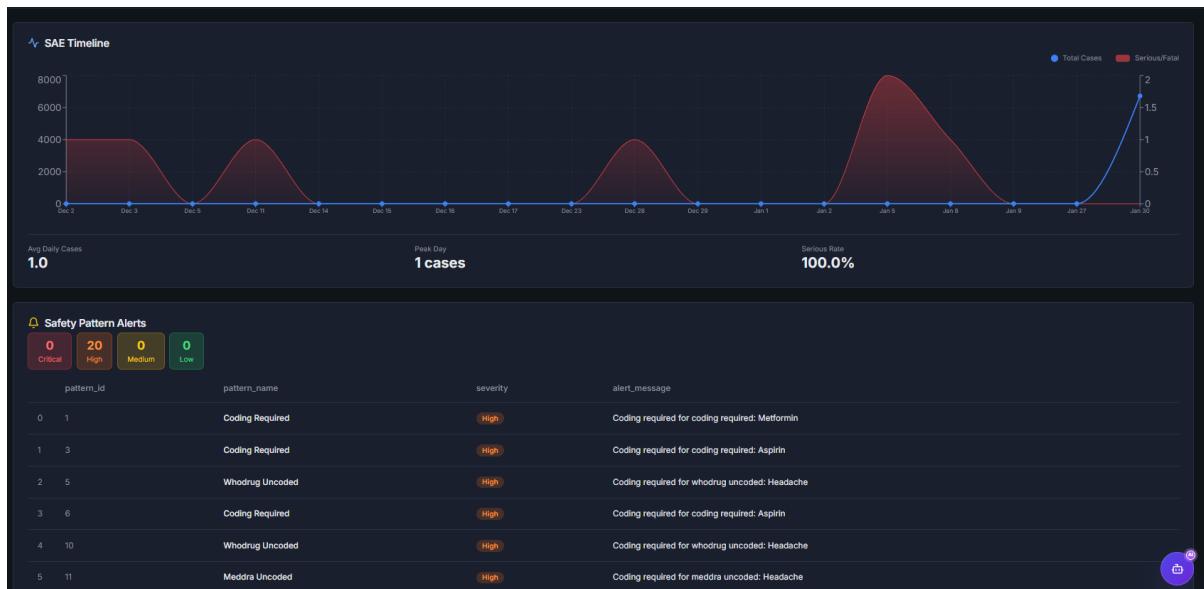


Figure 12: Safety Surveillance View 2: Timeline

6 Medical Coder View

6.1 Purpose

The Coder View streamlines medical coding workflows for MedDRA (adverse events) and WHO-Drug (medications) dictionary coding.

6.2 Key Features

- Coding Queue:** Prioritized list of terms awaiting coding

- **Dictionary Search:** Real-time search across MedDRA and WHODrug
- **Auto-Code Suggestions:** AI-suggested codes based on verbatim term
- **Approval Workflow:** One-click code acceptance or escalation
- **Productivity Metrics:** Coding volume and accuracy tracking
- **Uncoded Term Alerts:** Highlighting high-priority pending items

6.3 Coding Workflow

1. Review uncoded terms in priority queue
2. View AI-suggested codes with confidence scores
3. Accept suggestion or search for alternative code
4. Approve coding with electronic signature
5. System records audit trail and updates patient record

The screenshot shows the 'Medical Coder Workbench' interface. At the top, there's a header bar with the title 'Medical Coder Workbench' and a subtitle 'MedDRA and WHODrug coding with AI assistance'. A badge indicates '4311 Terms Pending'. Below the header are five summary cards: '4270 MedDRA Pending' (with a blue circular icon), '41 WHODrug Pending' (with a purple circular icon), '0 Coded Today' (with a green circular icon), '1724 High Conf. Ready' (with a yellow circular icon), and '749 Escalated' (with a red circular icon). A search bar at the bottom left says 'Search by verbatim term, patient ID, or site...'. Below the search bar is a navigation bar with tabs: 'Coding Queue (4311)', 'Escalations (749)', and 'Productivity'. The main content area is titled 'AI-Assisted Coding Queue' and contains a table of uncoded terms. The table has columns: Item ID, Dictionary, Verbatim Term, Suggested Code, Patient, Confidence, and Actions. Each row shows a list item with a checkbox, the dictionary name (MedDRA or WHODRUG), the verbatim term, a suggested code (e.g., Fever, Headache, Aspirin, Metformin, Aspirin, Nausea), the patient ID (PAT-XXXXXX), a confidence bar (e.g., 84%, 70%, 73%, 84%, 92%, 92%), and two buttons: 'Approve' and 'Reject'. There's also a 'Batch Approve High Confidence (10)' button at the top right of the table.

Figure 13: Medical Coder View 1: Overview

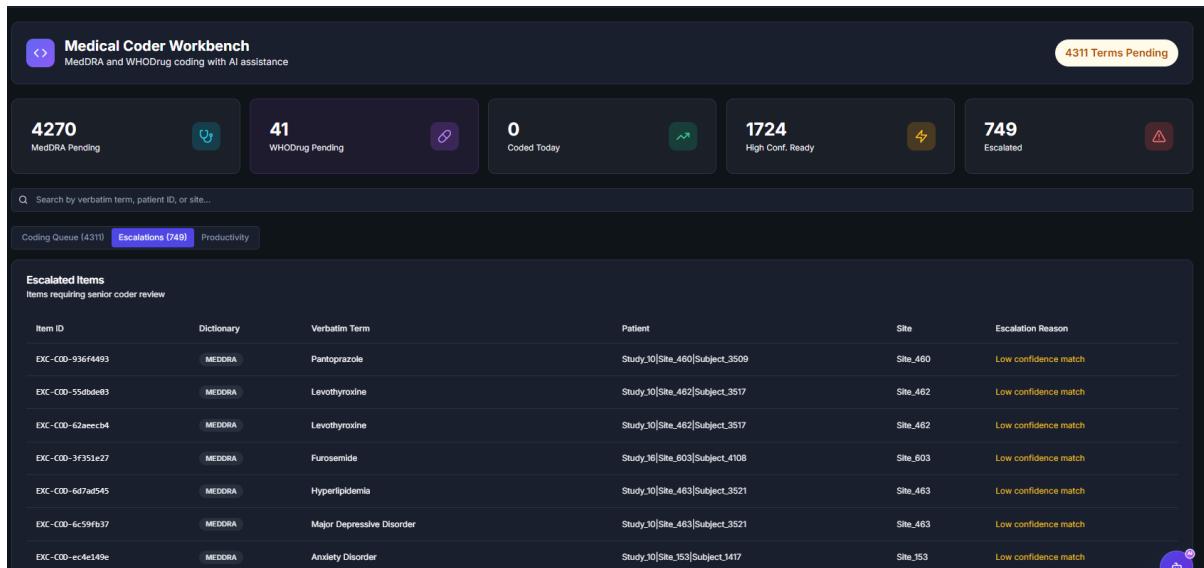


Figure 14: Medical Coder View 2: Escalations

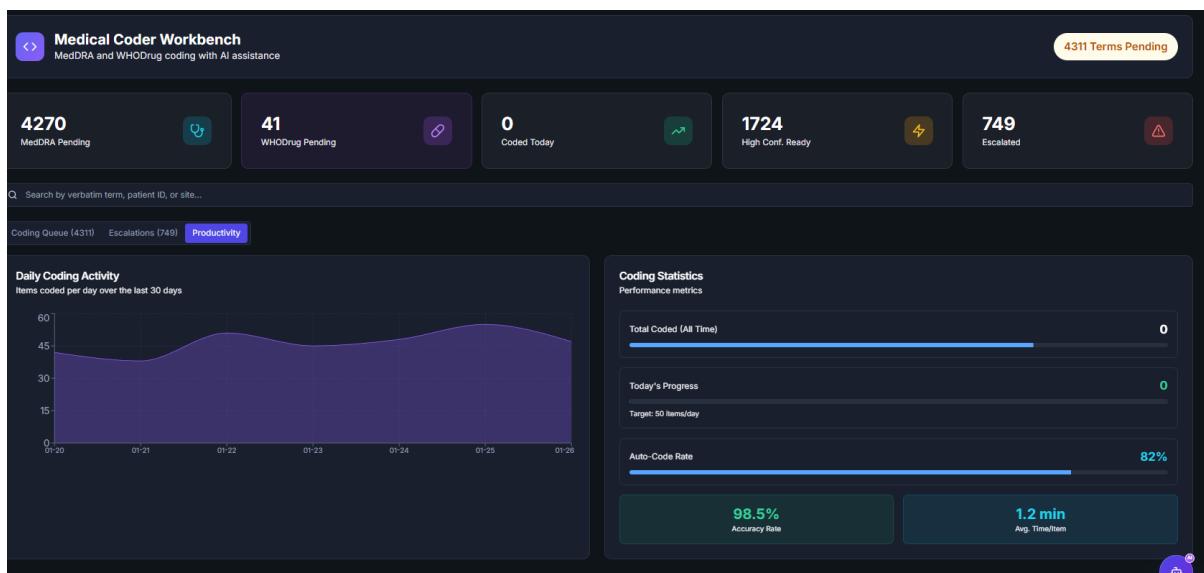


Figure 15: Medical Coder View 3: Statistics

7 Site Portal

7.1 Purpose

The Site Portal provides a site-centric view for site coordinators and local teams, focusing on their specific site's performance and patient population.

7.2 Key Features

- **Site Performance Summary:** DQI score, patient count, issue breakdown
- **Patient Registry:** Complete list of patients enrolled at the site

- **Issue Tracker:** Site-specific issues with resolution tools
- **Benchmark Comparison:** How the site compares to study average
- **Action Plan Generator:** AI-generated improvement recommendations
- **Activity Timeline:** Recent activity and audit history

7.3 Site Performance Indicators

Table 5: Site Portal KPIs

Indicator	Description
Site DQI Score	Aggregate data quality score (0-100)
Clean Patient Rate	Percentage of patients with no open issues
Query Resolution Time	Average days to resolve queries
Protocol Compliance	Adherence to visit windows and procedures

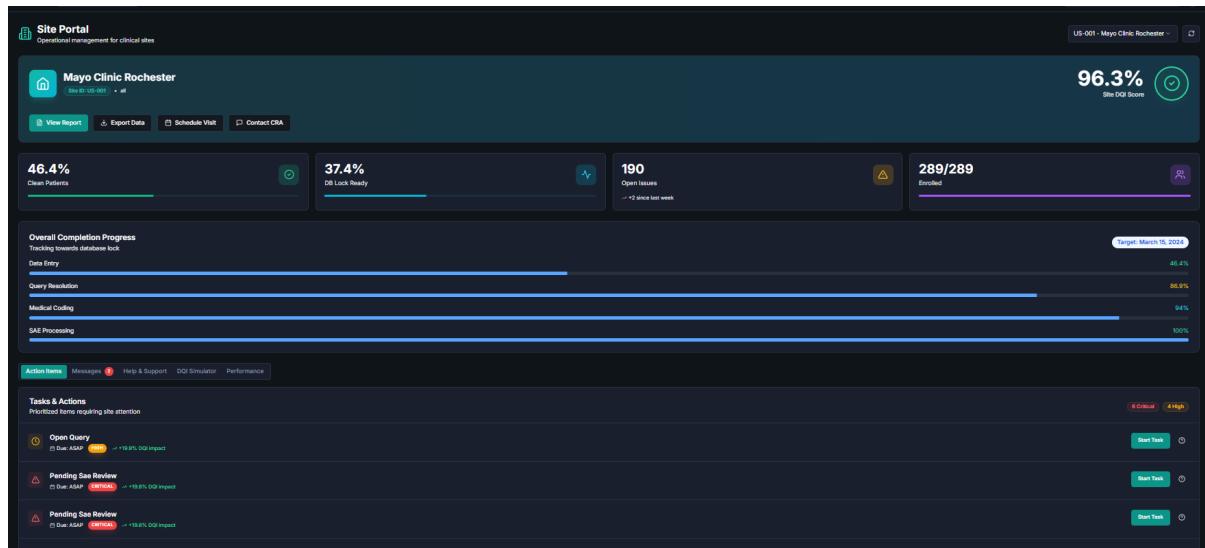


Figure 16: Site Portal View 1: Overview

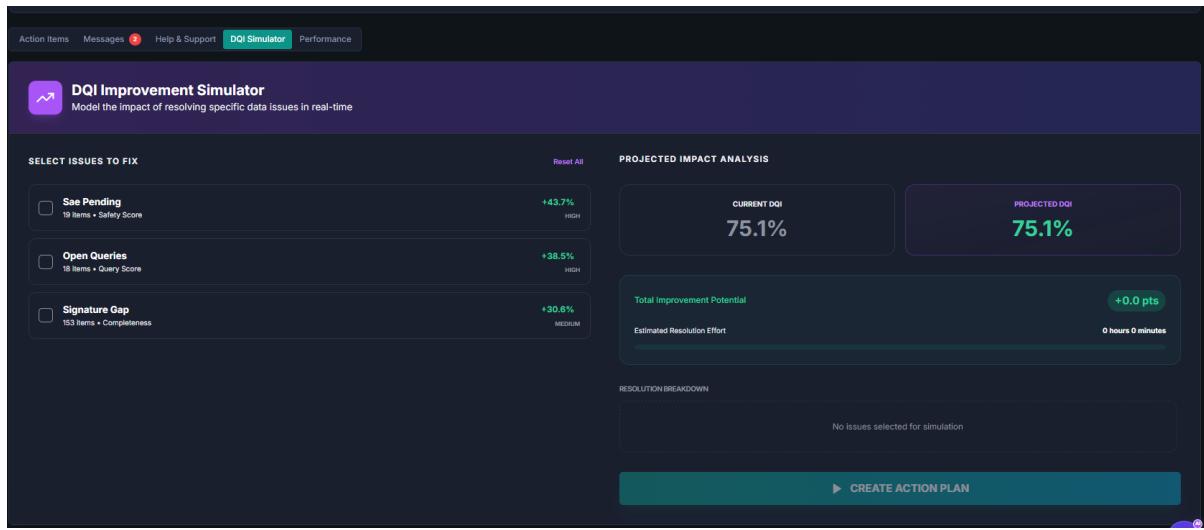


Figure 17: Site Portal View 2: DQI Improvement Simulation

8 Executive Overview

8.1 Purpose

The Executive Overview provides a high-level strategic view for senior leadership, focusing on portfolio health, major risks, and timeline projections.

8.2 Key Features

- **Portfolio Health Score:** Single metric summarizing overall status
- **Study Comparison Matrix:** Side-by-side study performance
- **Risk Heatmap:** Visual representation of risk distribution
- **Milestone Tracker:** Key milestone progress and projections
- **Regional Analysis:** Performance breakdown by geography
- **Trend Analysis:** Historical trends and forward projections



Figure 18: Executive View 1: Overview



Figure 19: Executive View 2: Regional, Site, DB Metrics

9 Collaboration Hub

9.1 Purpose

The Collaboration Hub enables real-time team communication and coordination across clinical operations roles, integrating messaging with clinical context.

9.2 Key Features

- Team Chat:** Role-based chat channels for clinical teams
- Issue Tagging:** Ability to reference specific patients or issues in messages

- **Activity Feed:** Real-time updates on issue resolutions and changes
- **Mention System:** @mention team members for attention
- **Thread Discussions:** Contextual conversations on specific topics
- **Notification Preferences:** Customizable alert settings per user

9.3 Communication Channels

The Collaboration Hub supports structured communication:

Table 6: Collaboration Channels

Channel Type	Purpose
Study Channel	Cross-functional study discussions
Site Channel	Site-specific coordination
Safety Channel	Pharmacovigilance team communication
Data Quality Channel	DM team coordination
General	Cross-study announcements

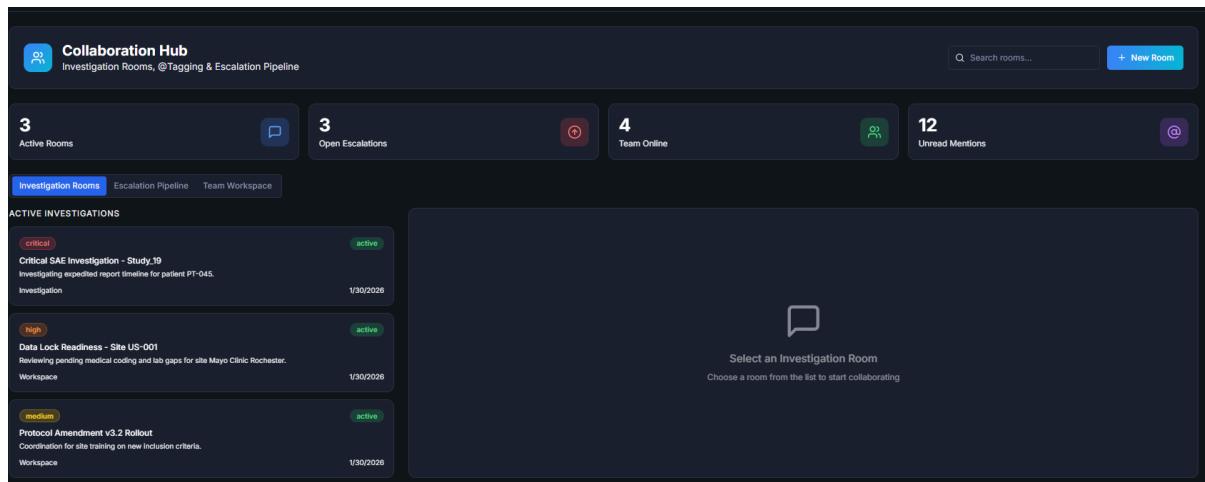


Figure 20: Collaboration Hub View 1: Overview

10 Cascade Explorer

10.1 Purpose

The Cascade Explorer visualizes the dependency graph between patients, issues, and sites, helping users understand the downstream impact of resolving specific issues.

10.2 Key Features

- **Interactive Graph:** Force-directed graph visualization of entities
- **Impact Analysis:** Shows how resolving one issue affects others
- **Filter Controls:** Filter by site, issue type, or patient status

- **Node Details:** Click any node to see detailed information
- **Cascade Path:** Trace blocking relationships between issues
- **Priority Highlighting:** Visual emphasis on high-impact nodes

10.3 Graph Visualization

The visualization uses a force-directed layout where:

- **Nodes** represent patients, issues, or sites
- **Edges** represent blocking relationships
- **Node Size** indicates impact score
- **Node Color** indicates risk level (red=critical, yellow=medium, green=low)

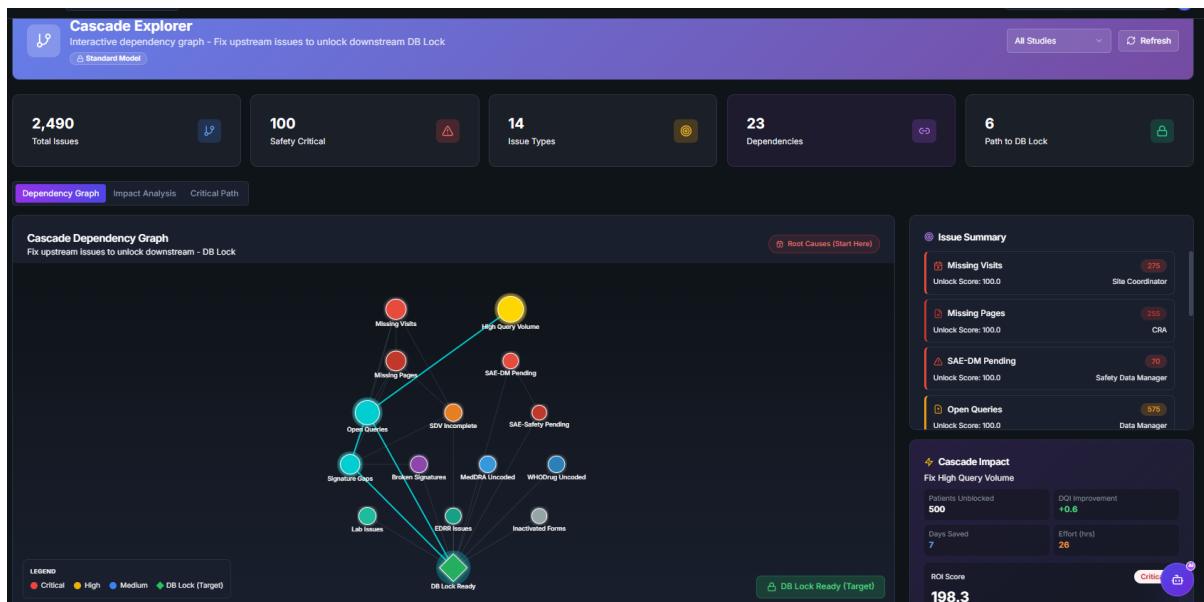


Figure 21: Cascade Explorer View 1: Overview

11 AI Assistant

11.1 Purpose

The AI Assistant provides a natural language interface for querying clinical trial data, available across all role-based views.

11.2 Key Features

- **Natural Language Queries:** Ask questions in plain English
- **Context Awareness:** Understands current study/site selection
- **Data Grounding:** Responses backed by live database queries
- **Agent Trace:** Shows which AI agents processed the query

- **Actionable Insights:** Recommendations with supporting evidence
- **Conversation Memory:** Maintains context within session

11.3 Example Queries

Users can ask questions such as:

- “Why is Site 05 underperforming?”
- “Which patients are blocking DB lock?”
- “Show me critical issues at Site 12”
- “What’s the projected DB lock date?”
- “Summarize safety concerns this week”

The AI Assistant routes queries through the 4-agent orchestrator (Supervisor → Diagnostic → Resolver → Communicator) to generate comprehensive responses.

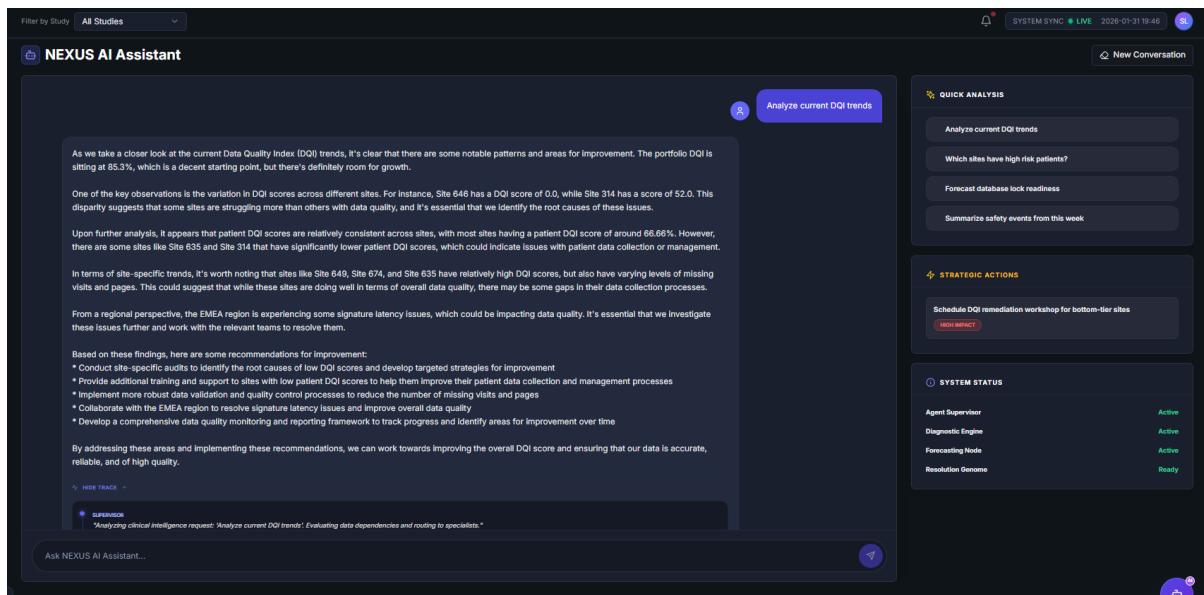


Figure 22: AI Assistant View 1: Overview

12 Reports & Export

12.1 Purpose

The Reports module enables users to generate and export professional reports in multiple formats for stakeholder communication.

12.2 Available Reports

Table 7: Report Types

Report	Description	Formats
Portfolio Summary	Overall portfolio status	PDF, DOCX
Site Performance	Site-level metrics	PDF, DOCX
Patient Status	Individual patient details	PDF
Risk Assessment	Risk tier breakdown	PDF, PPTX
DQI Analysis	Data quality deep-dive	PDF, DOCX
Executive Summary	Leadership briefing	PDF, PPTX

12.3 Report Generation Workflow

1. Select report type from available templates
2. Configure parameters (study, site, date range)
3. Preview report content in-browser
4. Generate final document in desired format
5. Download or share via configured channels

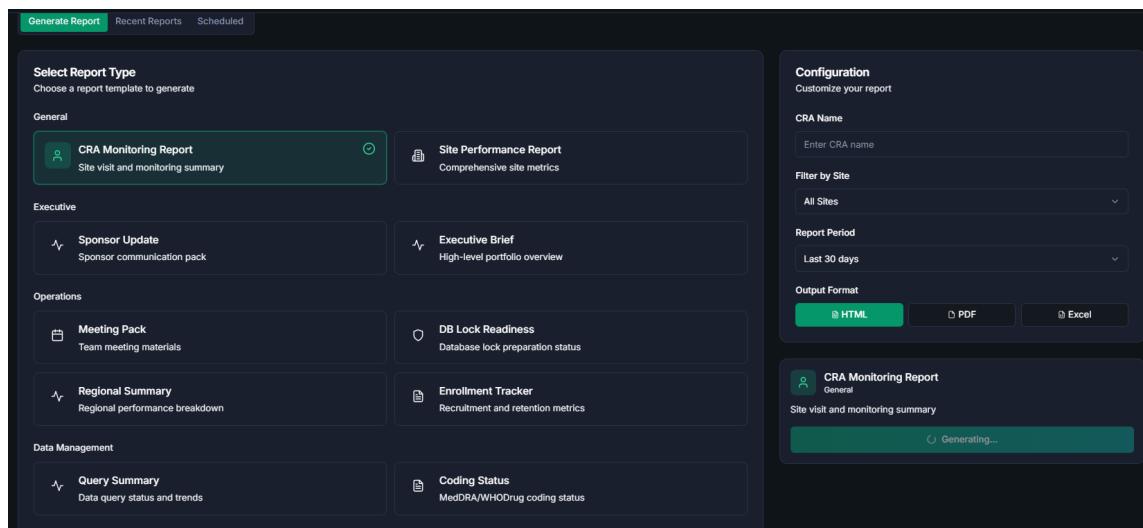


Figure 23: Reports - Generation Interface

13 ML Governance Dashboard

13.1 Purpose

The ML Governance Dashboard provides oversight for all machine learning models deployed in the system, ensuring regulatory compliance and model quality.

13.2 Key Features

- **Model Registry:** Complete list of all ML models with status

- **Approval Workflow:** Review and approve models before deployment
- **Drift Monitoring:** Detect when models need retraining
- **Performance Metrics:** Track model accuracy over time
- **Audit Trail:** Complete history of model changes
- **Deployment Controls:** Stage, deploy, or retire models

13.3 Model Lifecycle States

Table 8: Model Lifecycle

State	Description
Development	Model in training/testing
Pending Approval	Awaiting stakeholder review
Approved	Ready for deployment
Deployed	Active in production
Retired	Removed from service

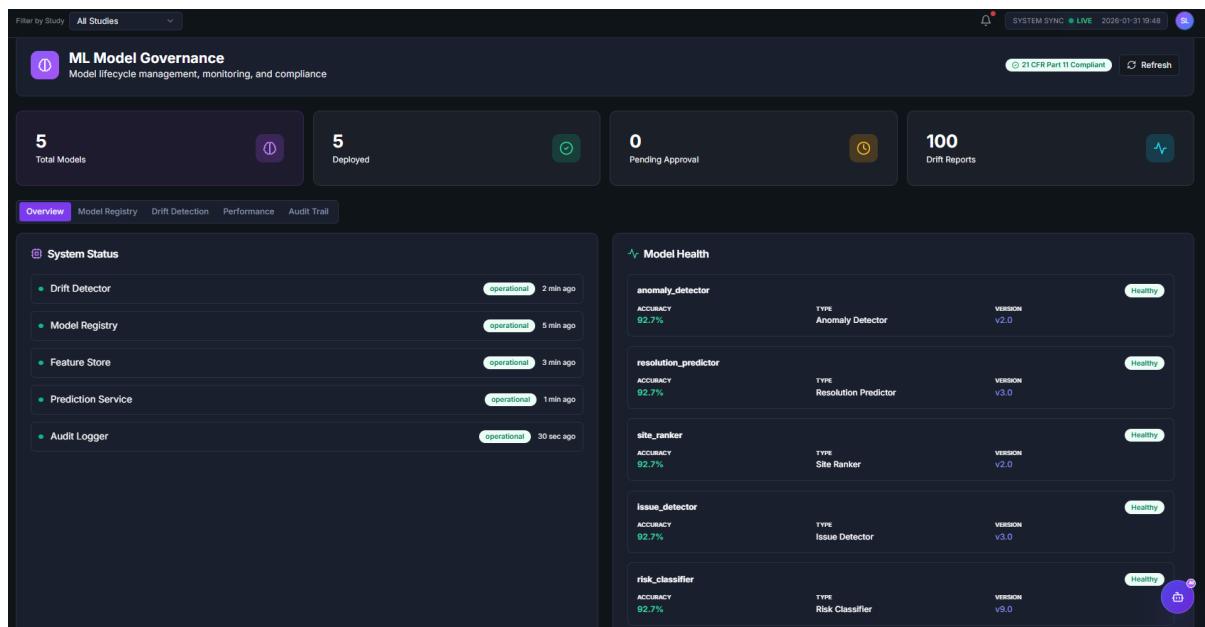


Figure 24: ML Governance View 1: Overview

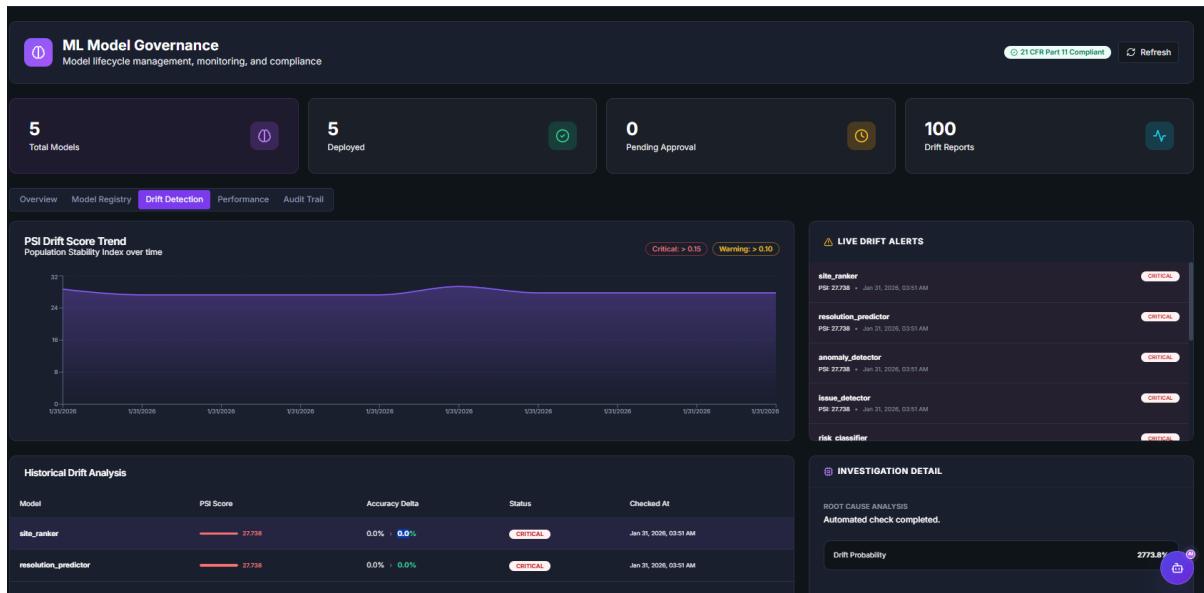


Figure 25: ML Governance View 2: PSI Drift Score & Live Report

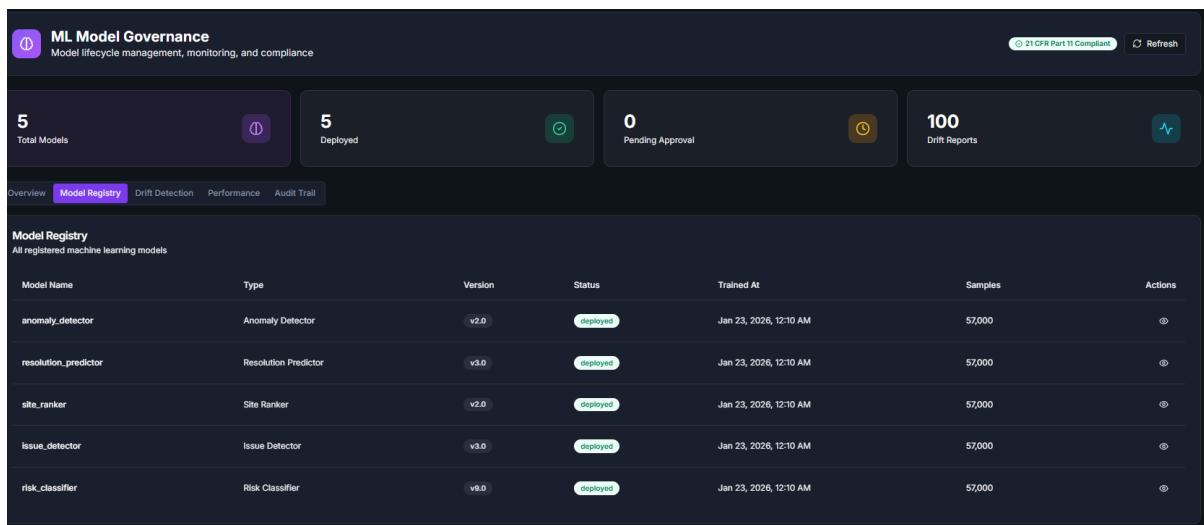


Figure 26: ML Governance View 3: Model Registry

14 Visualization Dashboard

14.1 Purpose

The Visualization Dashboard provides advanced analytics visualizations including charts, graphs, and interactive data exploration tools.

14.2 Visualization Types

- DQI Distribution:** Histogram of patient DQI scores
- Site Comparison:** Bar charts comparing site performance

- **Regional Heatmap:** Geographic performance visualization
- **Trend Lines:** Time-series charts for key metrics
- **Issue Breakdown:** Pie/donut charts by issue category
- **Forecast Projections:** DB lock and milestone projections

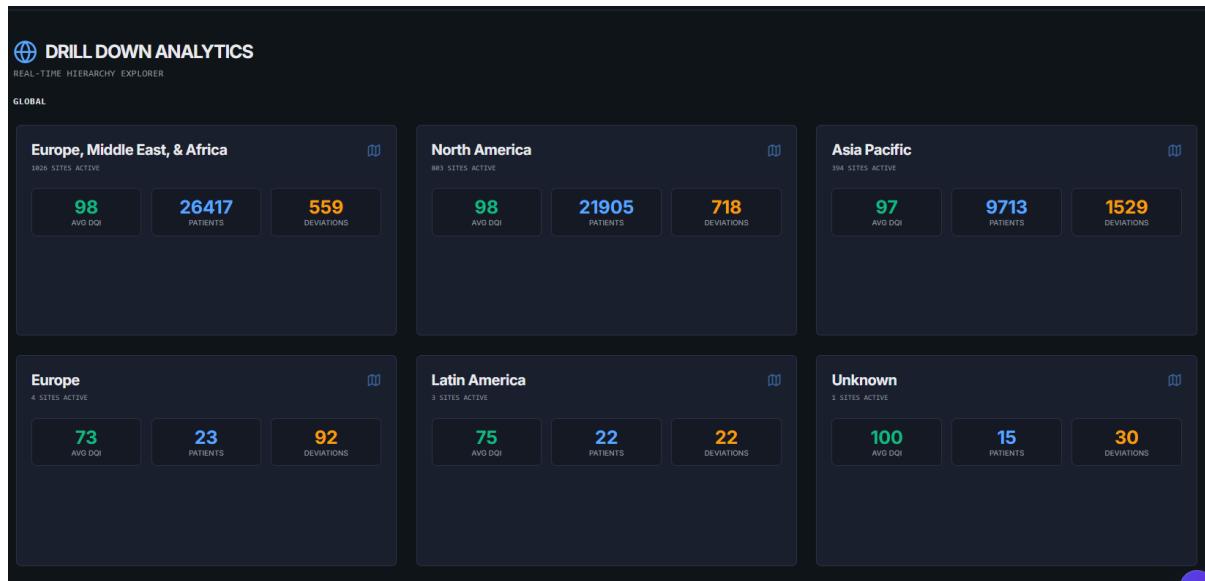


Figure 27: Drill Down Analytics View 1: Overview

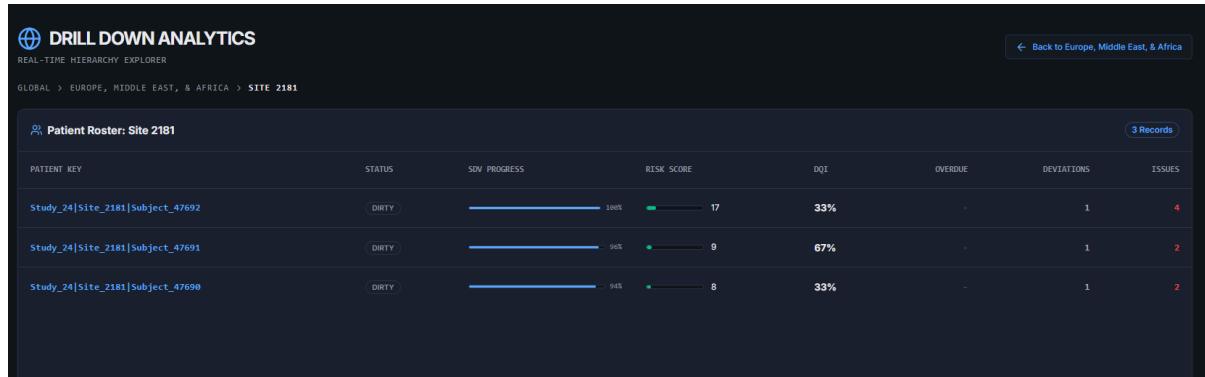


Figure 28: Drill Down Analytics View 2: Patients

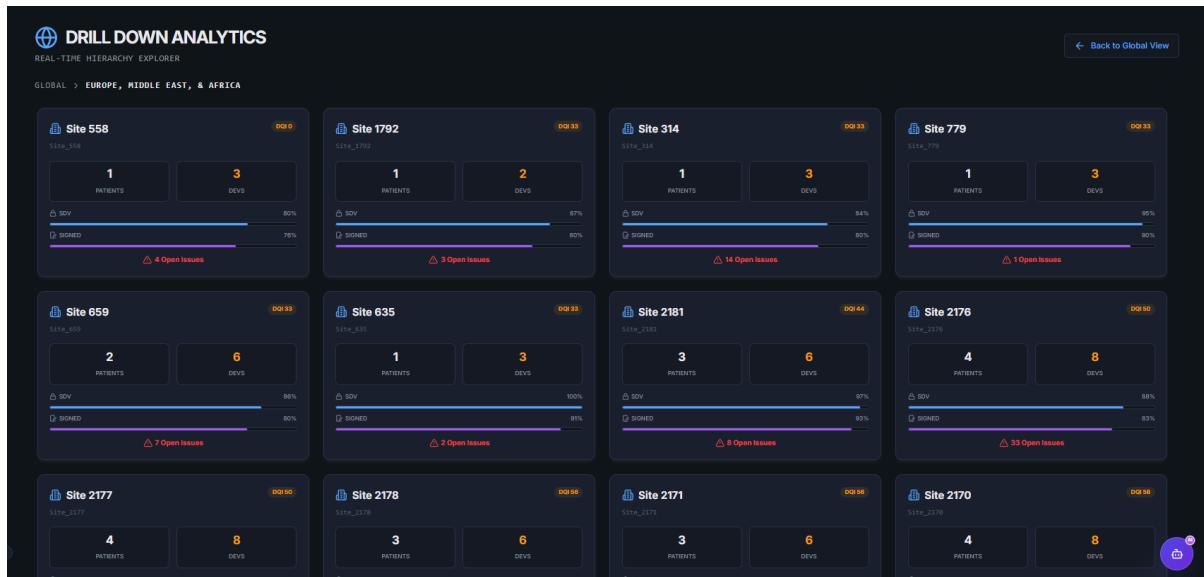


Figure 29: Drill Down Analytics View 3: Sites

TrialPlus Nexus 10X

Technical Documentation

Module 6: Operations & Reference

Scripts, Deployment, and API Reference

Version: 1.0.0
Date: January 2026
Scripts: 14 Pipeline Scripts
API Routes: 18 Route Modules

Contents

1 Pipeline Scripts	3
1.1 Overview	3
1.2 Data Pipeline Scripts	3
1.2.1 Data Ingestion	3
1.2.2 UPR Builder	3
1.3 Machine Learning Scripts	4
1.3.1 Training Workflow	4
1.4 Knowledge & Analytics Scripts	5
1.4.1 Embedding Pipeline	5
1.5 Utility Script	5
2 Deployment Guide	6
2.1 System Requirements	6
2.2 Environment Setup	6
2.2.1 Step 1: Clone Repository	6
2.2.2 Step 2: Python Environment	6
2.2.3 Step 3: Environment Variables	6
2.2.4 Step 4: Database Setup	6
2.2.5 Step 5: Frontend Dependencies	7
2.3 Running the Application	7
2.3.1 Unified Launcher (Recommended)	7
2.3.2 Manual Startup	7
2.4 Data Initialization	7
2.5 Accessing the Application	7
2.6 Default Credentials	7
2.7 Production Considerations	8
3 API Reference	8
3.1 Overview	8
3.2 Authentication	8
3.3 Patient Endpoints	9
3.4 Site Endpoints	9
3.5 Analytics Endpoints	9
3.6 Intelligence Endpoints	9
3.7 Issues Endpoints	10
3.8 Safety Endpoints	10
3.9 Coding Endpoints	10
3.10 Reports Endpoints	11
3.11 Simulation Endpoints	11
3.12 Graph Endpoints	11
3.13 Dashboard Endpoints	12
3.14 ML Endpoints	12
3.15 Agent Endpoints	12
3.16 Response Formats	12
3.17 Error Codes	13
3.18 Interactive Documentation	13

4 Future Features & Roadmap	13
4.1 Overview	13
4.2 Kafka Streaming Service	14
4.2.1 Purpose	14
4.2.2 Work Completed	14
4.2.3 Supported Event Types	14
4.2.4 Integration Path	14
4.3 Notification Service	14
4.3.1 Purpose	14
4.3.2 Work Completed	15
4.3.3 Notification Lifecycle	15
4.3.4 Integration Path	15
4.4 Electronic Signature Service	15
4.4.1 Purpose	15
4.4.2 Regulatory Compliance	16
4.4.3 Work Completed	16
4.4.4 Signature Meanings	16
4.4.5 Integration Path	17
4.5 LLM Fine-tuning Pipeline	17
4.5.1 Purpose	17
4.5.2 Work Completed	17
4.5.3 Training Data Categories	17
4.5.4 Integration Path	18
4.6 Roadmap Summary	18

1 Pipeline Scripts

1.1 Overview

The `src/` directory contains standalone Python scripts that execute specific pipeline stages. These scripts are designed to be run independently or as part of the unified launcher workflow.

1.2 Data Pipeline Scripts

Table 1: Data Pipeline Scripts

Script	Purpose
<code>run_ingestion.py</code>	Ingests raw clinical trial data from CSV files and loads into source tables (EDC, IVRS, Labs, Safety, Coding, Query, SDV, Signature)
<code>run_cleaning.py</code>	Applies data cleaning transformations including null handling, date parsing, and standardization
<code>run_upr_builder.py</code>	Builds the Unified Patient Record by joining all data sources into a single denormalized view per patient

1.2.1 Data Ingestion

The ingestion script reads from `data/raw/` and populates PostgreSQL tables. It handles 8 different clinical data domains and maintains referential integrity across studies and sites.

1.2.2 UPR Builder

The Unified Patient Record builder is the core ETL step that creates a 360-degree view of each patient by:

- Joining EDC data with visit schedules
- Merging safety events and lab results
- Aggregating query and SDV status
- Computing patient-level metrics

1.3 Machine Learning Scripts

Table 2: Machine Learning Training Scripts

Script	Purpose
<code>run_risk_training_v9_production.py</code>	Trains the patient risk classification model using XGBoost with 45 engineered features. Outputs to <code>models/risk_model_v9.pkl</code>
<code>run_resolution_time_PRODUCTION_v3.py</code>	Trains the ResolutionTimeModel that predicts days to resolve issues based on issue type and site characteristics
<code>run_anomaly_detector_OPTIMIZED_v2.py</code>	Trains Isolation Forest anomaly detection model to identify unusual patient or site patterns
<code>run_site_ranker_OPTIMIZED_v2.py</code>	Trains the SiteRiskRanker model that prioritizes sites needing attention
<code>run_issue_detector_production.py</code>	Trains issue classification model to categorize and prioritize data quality issues
<code>run_comprehensive_risk_training.py</code>	Full training pipeline that runs all ML models in sequence with cross-validation

1.3.1 Training Workflow

Each training script follows a consistent pattern:

1. Load processed data from `data/processed/` parquet files
2. Engineer features specific to the model's purpose
3. Split data with stratification for class balance
4. Train with hyperparameter tuning where applicable
5. Evaluate on held-out test set
6. Save model artifact with metadata to `models/`
7. Log metrics for drift monitoring

1.4 Knowledge & Analytics Scripts

Table 3: Knowledge Pipeline Scripts

Script	Purpose
<code>run_embedding_pipeline.py</code>	Generates vector embeddings for all text content using sentence-transformers. Populates ChromaDB collections
<code>run_rag_knowledge_base.py</code>	Indexes ICH-GCP guidelines, protocol documents, and SOPs for RAG retrieval
<code>run_causal_hypothesis.py</code>	Runs the Causal Hypothesis Engine to generate data-driven insights and patterns
<code>run_resolution_genome.py</code>	Analyzes historical resolution patterns to build the Resolution Genome knowledge base
<code>run_site_risk_ranker.py</code>	Generates site risk rankings for prioritized monitoring

1.4.1 Embedding Pipeline

The embedding pipeline processes 7 content categories:

- Patient issue descriptions
- Query text and resolutions
- Safety narratives
- Protocol sections
- ICH-GCP guidelines
- SOP procedures
- Historical resolution notes

Each category is stored in a separate ChromaDB collection with metadata for filtered retrieval.

1.5 Utility Script

Table 4: Utility Script

Script	Purpose
<code>rebuild_upr_perfect.py</code>	Complete rebuild utility that regenerates the entire Unified Patient Record from scratch, useful after schema changes

2 Deployment Guide

2.1 System Requirements

Table 5: Minimum System Requirements

Component	Requirement
Operating System	Windows 10/11, macOS, or Linux
Python	3.10+
Node.js	18+
PostgreSQL	14+
Neo4j	5.x (optional, for graph features)
RAM	8GB minimum, 16GB recommended
Disk Space	10GB for application + data

2.2 Environment Setup

2.2.1 Step 1: Clone Repository

Clone the TrialPlus Nexus repository from the version control system to your local machine.

2.2.2 Step 2: Python Environment

Create a virtual environment and install dependencies:

- Create virtual environment in project root
- Activate the environment
- Install requirements from `requirements.txt`

2.2.3 Step 3: Environment Variables

Create a `.env` file in the project root with the following configuration:

Table 6: Required Environment Variables

Variable	Description
<code>DATABASE_URL</code>	PostgreSQL connection string
<code>GOOGLE_API_KEY</code>	Google Gemini API key for LLM
<code>NEO4J_URI</code>	Neo4j connection URI (optional)
<code>NEO4J_USER</code>	Neo4j username (optional)
<code>NEO4J_PASSWORD</code>	Neo4j password (optional)
<code>JWT_SECRET_KEY</code>	Secret for JWT token signing

2.2.4 Step 4: Database Setup

Initialize the PostgreSQL database:

1. Create a new PostgreSQL database
2. Run the setup script which creates all required tables
3. The script automatically creates 25+ tables with proper indexes

2.2.5 Step 5: Frontend Dependencies

Navigate to the frontend directory and install Node.js dependencies using npm.

2.3 Running the Application

2.3.1 Unified Launcher (Recommended)

The project includes a unified launcher script that starts all components:

Table 7: Unified Launcher Commands

Command	Action
Run launcher	Starts backend (port 8000) and frontend (port 5173)
Backend only	Starts only the FastAPI backend
Frontend only	Starts only the React frontend

2.3.2 Manual Startup

To start components manually:

1. **Backend:** Navigate to project root and run uvicorn with the FastAPI app
2. **Frontend:** Navigate to frontend directory and run npm dev server

2.4 Data Initialization

After starting the application for the first time:

1. **Load Raw Data:** Place CSV files in `data/raw/` following the expected schema
2. **Run Ingestion:** Execute the ingestion script to populate source tables
3. **Build UPR:** Run the UPR builder to create unified patient records
4. **Train Models:** Run ML training scripts or use pre-trained models from `models/`
5. **Initialize Knowledge:** Run embedding and RAG scripts for AI features

2.5 Accessing the Application

Table 8: Application URLs

Component	URL	Purpose
Frontend	http://localhost:5173	React application
Backend API	http://localhost:8000	FastAPI REST API
API Docs	http://localhost:8000/docs	Swagger UI documentation
ReDoc	http://localhost:8000/redoc	Alternative API docs

2.6 Default Credentials

For demo purposes, the following test users are available:

Table 9: Demo User Accounts

Username	Password	Role	Access
admin	admin123	Administrator	Full access + admin
lead	lead123	Study Lead	Portfolio view, approvals
dm	dm123	Data Manager	Query management, data editing
cra	cra123	CRA	Site monitoring, reports
coder	coder123	Medical Coder	MedDRA/WHO Drug coding
safety	safety123	Safety Officer	SAE tracking, narratives
exec	exec123	Executive	Dashboard views, reports
testuser	testpassword	Test User	General testing

Note: In TEST_MODE, the system accepts any username/password combination and creates users on-the-fly. For production, use the exact credentials listed above.

2.7 Production Considerations

For production deployment:

- **Security:** Change all default passwords and use strong JWT secrets
- **HTTPS:** Configure SSL/TLS certificates for encrypted connections
- **Database:** Use connection pooling and configure proper backups
- **Monitoring:** Set up application monitoring and log aggregation
- **Scaling:** Consider horizontal scaling with load balancer for high traffic

3 API Reference

3.1 Overview

The TrialPlus Nexus API is built with FastAPI and provides a comprehensive REST interface for all system functionality. The API is organized into 18 route modules covering different functional domains.

3.2 Authentication

All API endpoints (except login) require JWT Bearer token authentication. Tokens are obtained via the login endpoint and must be included in the Authorization header.

Table 10: Authentication Endpoints

Method	Endpoint	Description
POST	/api/v1/auth/login	Authenticate and receive JWT tokens
POST	/api/v1/auth/refresh	Refresh expired access token
POST	/api/v1/auth/logout	Invalidate current session
GET	/api/v1/auth/me	Get current user profile

3.3 Patient Endpoints

The patients API provides access to patient data, issues, and risk information.

Table 11: Patient API Endpoints

Method	Endpoint	Description
GET	/patients	List patients with pagination
GET	/patients/{key}	Get single patient details
GET	/patients/search	Search patients by query
GET	/patients/dqi	Get DQI distribution
GET	/patients/clean-status	Get clean patient counts
GET	/patients/dblock-status	Get DB lock readiness
GET	/patients/issues	Get patient issue summary
GET	/patients/{key}/risk-explanation	Get ML risk explanation

3.4 Site Endpoints

Site management and monitoring endpoints.

Table 12: Site API Endpoints

Method	Endpoint	Description
GET	/sites	List all sites
GET	/sites/{id}	Get site details
GET	/sites/benchmarks	Get site benchmarks
GET	/sites/smart-queue	AI-prioritized work queue
GET	/sites/{id}/patients	Get patients at site
GET	/sites/{id}/portal	Get site portal data
GET	/sites/{id}/dqi-issues	Get site DQI issues
POST	/sites/{id}/action-plan	Generate action plan

3.5 Analytics Endpoints

Data analytics and visualization endpoints.

Table 13: Analytics API Endpoints

Method	Endpoint	Description
GET	/analytics/portfolio	Portfolio summary metrics
GET	/analytics/dqi-distribution	DQI score distribution
GET	/analytics/cascade	Cascade analysis data
GET	/analytics/patterns	Detected patterns
GET	/analytics/bottlenecks	Process bottlenecks
GET	/analytics/quality-matrix	Quality matrix data
GET	/analytics/recommendations	AI recommendations
POST	/analytics/refresh	Refresh analytics cache

3.6 Intelligence Endpoints

AI and agentic intelligence endpoints.

Table 14: Intelligence API Endpoints

Method	Endpoint	Description
GET	/intelligence/hypotheses	Get causal hypotheses
POST	/intelligence/swarm/run	Run agent swarm investigation
GET	/intelligence/anomalies	Get detected anomalies
POST	/intelligence/auto-fix	Auto-resolve an issue
POST	/intelligence/assistant/query	Natural language query
POST	/intelligence/assistant/reset	Reset AI session

3.7 Issues Endpoints

Issue tracking and resolution endpoints.

Table 15: Issues API Endpoints

Method	Endpoint	Description
GET	/issues	List issues with filters
GET	/issues/summary	Issue summary counts
POST	/issues	Create new issue
PUT	/issues/{id}	Update issue
POST	/issues/{id}/resolve	Resolve issue
POST	/issues/{id}/escalate	Escalate issue

3.8 Safety Endpoints

Pharmacovigilance and SAE management endpoints.

Table 16: Safety API Endpoints

Method	Endpoint	Description
GET	/safety/overview	Safety dashboard data
GET	/safety/sae-cases	List SAE cases
GET	/safety/sla-status	Regulatory SLA status
GET	/safety/signals	Safety signal detection
GET	/safety/timeline	SAE timeline data
GET	/safety/narratives/{id}	Get SAE narrative
POST	/safety/sae/{id}/update-status	Update SAE status

3.9 Coding Endpoints

Medical coding (MedDRA/WHODrug) endpoints.

Table 17: Coding API Endpoints

Method	Endpoint	Description
GET	/coding/queue	Get coding queue
GET	/coding/meddra/pending	MedDRA pending items
GET	/coding/whodrug/pending	WHODrug pending items
GET	/coding/stats	Coding statistics
POST	/coding/approve/{id}	Approve coding
GET	/coding/search/{dict}	Search dictionary

3.10 Reports Endpoints

Report generation endpoints.

Table 18: Reports API Endpoints

Method	Endpoint	Description
GET	/reports/types	List available report types
POST	/reports/generate	Generate a report
GET	/reports/preview/{type}	Preview report content

3.11 Simulation Endpoints

Digital twin and what-if simulation endpoints.

Table 19: Simulation API Endpoints

Method	Endpoint	Description
GET	/simulation/scenarios	List scenario types
POST	/simulation/run	Run simulation
POST	/simulation/compare	Compare scenarios
GET	/simulation/current-state	Get current state
GET	/simulation/projections	Get projections
GET	/simulation/db-lock-projection	DB lock timeline
POST	/simulation/what-if	Run what-if analysis

3.12 Graph Endpoints

Neo4j knowledge graph and cascade endpoints.

Table 20: Graph API Endpoints

Method	Endpoint	Description
GET	/graph/nodes	Get graph nodes
GET	/graph/edges	Get graph edges
GET	/graph/cascade-path/{id}	Get cascade path
GET	/graph/cascade-analysis	Full cascade analysis
GET	/graph/visualization-data	Data for graph viz

3.13 Dashboard Endpoints

Role-based dashboard data endpoints.

Table 21: Dashboard API Endpoints

Method	Endpoint	Description
GET	/dashboards/summary	General summary
GET	/dashboards/main	Main dashboard
GET	/dashboards/cra	CRA dashboard data
GET	/dashboards/data_manager	DM dashboard data
GET	/dashboards/safety	Safety dashboard data
GET	/dashboards/study_lead	Study Lead data
GET	/dashboards/site	Site dashboard data
GET	/dashboards/coder	Coder dashboard data

3.14 ML Endpoints

Machine learning model management endpoints.

Table 22: ML API Endpoints

Method	Endpoint	Description
GET	/ml/models	List ML models
GET	/ml/summary	ML summary stats
POST	/ml/models/{id}/approve	Approve model
POST	/ml/models/{id}/deploy	Deploy model
GET	/ml/drift-reports	Get drift reports
GET	/ml/audit-log	Model audit log

3.15 Agent Endpoints

Direct agent invocation endpoints.

Table 23: Agents API Endpoints

Method	Endpoint	Description
GET	/agents	List all agents
GET	/agents/supervisor/status	Supervisor status
POST	/agents/supervisor/act	Invoke supervisor
POST	/agents/diagnostic/act	Invoke diagnostic
POST	/agents/forecaster/act	Invoke forecaster
POST	/agents/resolver/act	Invoke resolver
POST	/agents/executor/act	Invoke executor
POST	/agents/communicator/act	Invoke communicator

3.16 Response Formats

All API responses follow a consistent JSON structure:

Table 24: Standard Response Fields

Field	Type	Description
success	boolean	Whether request succeeded
data	object/array	Response payload
message	string	Human-readable message
error	string	Error description if failed
timestamp	string	ISO 8601 timestamp

3.17 Error Codes

Table 25: HTTP Status Codes

Code	Meaning
200	Success
201	Created successfully
400	Bad request (invalid parameters)
401	Unauthorized (missing/invalid token)
403	Forbidden (insufficient permissions)
404	Resource not found
422	Validation error
500	Internal server error

3.18 Interactive Documentation

FastAPI automatically generates interactive API documentation accessible at:

- **Swagger UI:** `http://localhost:8000/docs` – Interactive API explorer with try-it-out functionality
- **ReDoc:** `http://localhost:8000/redoc` – Clean, readable API reference

4 Future Features & Roadmap

4.1 Overview

The following modules have been developed and exist in the codebase but are not yet integrated into the running demo. These represent infrastructure investments for future capabilities and production readiness.

Table 26: Future Feature Modules

Module	Directory	Lines of Code	Status
Kafka Streaming	<code>src/streaming/</code>	900+	Ready for integration
Notification Service	<code>src/notifications/</code>	660+	Ready for integration
Electronic Signatures	<code>src/compliance/</code>	480+	Ready for integration
LLM Fine-tuning	<code>src/finetuning/</code>	400+	Development

4.2 Kafka Streaming Service

4.2.1 Purpose

The streaming module provides real-time event streaming infrastructure using Apache Kafka, enabling live updates to dashboards and triggering automated workflows when data changes occur.

4.2.2 Work Completed

The streaming service includes the following components:

- **Kafka Producer:** Publishes events to Kafka topics with delivery guarantees. Supports batch sending and compression for high throughput.
- **Kafka Consumer:** Subscribes to topics and processes incoming events with at-least-once semantics. Supports consumer groups for parallel processing.
- **Event Processor:** Routes incoming events to appropriate handlers based on event type. Implements dead letter queue for failed messages.
- **UPR Updater:** Listens for data change events and incrementally updates the Unified Patient Record without full rebuild.
- **Stream Configuration:** Centralized configuration for topics, partitions, and retention policies.

4.2.3 Supported Event Types

Table 27: Streaming Event Types

Event Type	Description
Patient Update	Patient record field changes
Issue Detected	New issue identified by ML models
Issue Resolved	Issue marked as resolved
DQI Change	Patient DQI score recalculated
Site Update	Site-level metric changes
Drift Detected	ML model drift alert

4.2.4 Integration Path

To enable streaming in production:

1. Deploy Apache Kafka cluster (or use managed service like Confluent Cloud)
2. Configure Kafka broker addresses in environment variables
3. Enable streaming module in backend configuration
4. Update frontend to use WebSocket connections for real-time updates

4.3 Notification Service

4.3.1 Purpose

The notification service provides a multi-channel communication system for alerting users about important events, including in-app notifications, email, SMS, and integration with collaboration platforms like Slack and Microsoft Teams.

4.3.2 Work Completed

The notification module includes:

- **Notification Queue:** Priority-based queue with PostgreSQL persistence. Supports urgent, high, normal, and low priority levels.
- **Multi-Channel Delivery:** Abstracted delivery interface supporting in-app, email, SMS, push notifications, Slack, and Teams.
- **Notification Types:** Predefined types including alerts, mentions, escalations, tasks, issues, safety events, system messages, digests, and reports.
- **User Preferences:** Per-user configuration for preferred channels, quiet hours, batch preferences, and rate limiting.
- **Delivery Tracking:** Full lifecycle tracking from pending to delivered to read with timestamps.
- **Batch Digests:** Aggregation of multiple notifications into daily or weekly digest emails.

4.3.3 Notification Lifecycle

Notifications progress through the following states:

Table 28: Notification Status Lifecycle

Status	Description
Pending	Created, not yet processed
Queued	Added to delivery queue
Sent	Dispatched to channel
Delivered	Confirmed delivery
Read	User acknowledged
Archived	Moved to archive
Failed	Delivery failed

4.3.4 Integration Path

To enable notifications in production:

1. Configure SMTP server for email delivery
2. Set up SMS gateway (Twilio, AWS SNS)
3. Create Slack/Teams apps and configure webhooks
4. Enable notification service in backend
5. Add notification UI components to frontend

4.4 Electronic Signature Service

4.4.1 Purpose

The electronic signature module implements 21 CFR Part 11 compliant digital signatures for regulatory submissions. This is essential for GxP-regulated clinical trial environments where electronic records require legally binding signatures.

4.4.2 Regulatory Compliance

The service implements the following 21 CFR Part 11 requirements:

Table 29: 21 CFR Part 11 Compliance Features

Requirement	Implementation
User Identity Verification	Password authentication plus multi-factor authentication (MFA) via TOTP, SMS, or hardware token
Signature Meaning	Predefined meanings: Authored, Reviewed, Approved, Verified, Witnessed, Acknowledged, Confirmed
Trusted Timestamp	Server-side timestamps with time zone and synchronization
Immutable Hash	SHA-256 HMAC signature hash with chained integrity
Audit Trail	Complete record of all signature events with tamper detection

4.4.3 Work Completed

The electronic signature service includes:

- **Signature Creation:** Full workflow for creating legally binding e-signatures with MFA verification
- **Hash Chain Integrity:** Each signature references the previous signature hash, creating an immutable chain that detects tampering
- **Verification Service:** Validates signature authenticity by recomputing hash and checking chain integrity
- **Audit Trail:** Complete history of signature events with checksums for tamper detection
- **Multi-Record Support:** Sign any record type (issues, reports, approvals) with consistent interface

4.4.4 Signature Meanings

Per 21 CFR Part 11, each signature includes a legally binding meaning:

Table 30: Standard Signature Meanings

Meaning	Legal Statement
Authored	“I authored this record”
Reviewed	“I reviewed this record”
Approved	“I approve this record”
Verified	“I verified the accuracy”
Witnessed	“I witnessed this action”

4.4.5 Integration Path

To enable electronic signatures:

1. Implement MFA integration (Google Authenticator, Authy, or hardware tokens)
2. Create signature capture UI with meaning selection
3. Add signature verification endpoints to API
4. Integrate signature requirements into approval workflows
5. Configure audit trail export for regulatory inspections

4.5 LLM Fine-tuning Pipeline

4.5.1 Purpose

The fine-tuning module provides infrastructure for training custom language models on clinical trial data, enabling more accurate and domain-specific AI responses.

4.5.2 Work Completed

The fine-tuning pipeline includes:

- **Training Data Generator:** Automatically generates training examples from the PostgreSQL database covering DQI analysis, issue resolution, site performance, and general queries
- **Data Preparation:** Converts clinical data into instruction-following format compatible with modern LLM fine-tuning (Alpaca and chat formats)
- **Configuration Management:** Centralized configuration for model parameters, training hyperparameters, and system prompts
- **Jupyter Notebook:** Interactive notebook for Google Colab-based training using parameter-efficient fine-tuning (LoRA/QLoRA)
- **Training Scripts:** Command-line scripts for running fine-tuning jobs on local or cloud infrastructure

4.5.3 Training Data Categories

The generator creates examples across five categories:

Table 31: Fine-tuning Data Categories

Category	Description
DQI Analysis	Explaining patient data quality scores and contributing factors
Issue Resolution	Recommending resolution steps for different issue types
Site Performance	Analyzing site metrics and providing improvement suggestions
Query Response	Answering natural language questions about clinical data
General Analysis	Broad analytical reasoning about trial performance

4.5.4 Integration Path

To use fine-tuned models:

1. Generate training data using the data preparation script
2. Upload training data to Google Colab or training infrastructure
3. Run fine-tuning notebook with appropriate base model (Gemma, Llama, etc.)
4. Export fine-tuned model weights
5. Update LLM wrapper to use fine-tuned model for inference
6. Compare performance against base model on held-out test set

4.6 Roadmap Summary

Table 32: Feature Integration Priority

Feature	Priority	Effort	Value
Notification Service	High	Medium	User engagement, timely alerts
Electronic Signatures	High	Medium	Regulatory compliance
Kafka Streaming	Medium	High	Real-time updates
LLM Fine-tuning	Low	High	Improved AI accuracy