

Contents

Overview	2
Key Features	2
Dataset Visualization Overview	2
System Architecture	2
Scenario Coverage & Distribution	3
Scenario Breakdown	3
Dataset Generation Workflow	3
Research Impact & Applications	3
Dataset Specifications	4
Core Statistics	4
Technical Coverage	4
Environmental Conditions	4
Parameter Categories	4
Atmospheric Turbulence (5 parameters)	4
Weather & Environment (4 parameters)	4
Fading & Channel (4 parameters)	4
Pointing & Alignment (4 parameters)	5
OAM Beam Physics (4 parameters)	5
Hardware Impairments (3 parameters)	5
Propagation Physics (4 parameters)	5
Performance Metrics (5 parameters)	5
File Structure	5
Quick Start	6
Load Dataset	6
Basic Analysis	6
Visualization	6
Validation Results	7
Quality Metrics	7
Applications	7
Deep Reinforcement Learning	7
6G System Design	7
Research Applications	8
Industry Applications	8
Competition Instructions and Challenges	8
Overview	8
Competition Framework	8
Submission Requirements	11
Analysis Tools and Utilities	12
Competition Timeline and Prizes	14
Getting Started	14
Standards Compliance	15
6G Requirements	15
ITU-R IMT-2030 Framework	15
Citation	15
License & Usage	16
Usage Rights	16

Attribution Requirements	16
Contact	16

6G OAM THz Dataset

Dataset	270K samples	Physics Models	33 parameters	Quality Score	1.0/1.0	Frequency	300-600 GHz
---------	--------------	----------------	---------------	---------------	---------	-----------	-------------

Overview

The **6G OAM THz Dataset** is the world’s first comprehensive, physics-based dataset specifically designed for **Orbital Angular Momentum (OAM) beams in Terahertz (THz) frequencies** for Deep Reinforcement Learning applications in 6G wireless communications.

Key Features

- **270,000 high-fidelity samples** with perfect validation quality (1.0/1.0)
- **33 comprehensive physics parameters** covering complete THz propagation
- **4 deployment scenarios** from lab to wide-area coverage
- **300-600 GHz frequency range** with realistic atmospheric modeling
- **ITU-R IMT-2030 compliant** meeting 6G requirements

Dataset Visualization Overview

Dataset Comprehensive Overview

Figure 1: Dataset Comprehensive Overview

Figure 1: Complete Dataset Parameter Coverage - Comprehensive visualization showing all 33 physics parameters organized across 8 domains: Atmospheric Turbulence, Weather Environment, Fading & Channel, Pointing & Alignment, OAM Beam Physics, Hardware Impairments, Propagation Physics, and Performance Metrics. Each heatmap panel demonstrates normalized parameter distributions across 270K samples.

System Architecture

System Architecture

Figure 2: System Architecture

Figure 2: 6G OAM THz System Architecture - Technical architecture showing the complete system design with four main layers: Input Layer (environmental and hardware parameters), Physics Engine Core (atmospheric models, OAM beam physics, channel modeling), Hardware Layer (RF components and signal processing), and Propagation Channel (THz wave propagation with realistic impairments).

Scenario Coverage Matrix

Figure 3: Scenario Coverage Matrix

Scenario Coverage & Distribution

Figure 3: Comprehensive Scenario Coverage Matrix - Detailed breakdown of all 270,000 samples across 9 deployment scenarios with exact sample counts, parameter ranges, SINR coverage, throughput capabilities, latency performance, and validation scores for each scenario type.

Scenario Breakdown

Scenario	Sample Count	Percentage	Distance Range	Primary Use Case
Lab Controlled	70,000	25.9%	1-100m	Algorithm development, baseline validation
Indoor Realistic	80,000	29.6%	1-200m	Enterprise networks, smart buildings
Outdoor Urban	70,000	25.9%	10-1000m	Urban hotspots, backhaul connections
High Mobility	50,000	18.5%	50-5000m	Vehicle communications, high-speed scenarios

Dataset Generation Workflow

Dataset Workflow

Figure 4: Dataset Workflow

Figure 4: Dataset Generation and Validation Pipeline - Complete workflow showing the three-phase process: Configuration Phase (YAML setup and parameter validation), Generation Loop (physics simulation, scenario sampling, and quality checks), and Validation Pipeline (consistency analysis, range verification, and final scoring).

Research Impact & Applications

Research Impact

Figure 5: Research Impact

Figure 5: Global Research Impact and Applications - Comprehensive impact diagram showing six main research branches: Deep Reinforcement Learning (32-state action space, 500 steps exploration), 6G System Optimization (beam steering, power allocation), Standards Compliance Research (6G compliance, ITU-R alignment), Physics-Based Channel Modeling (atmospheric physics, OAM beam physics), Academic Publications (50+ paper potential), and Industry Testbed Development (equipment validation).

Dataset Specifications

Core Statistics

- **Total Samples:** 270,000
- **Parameters:** 33 physics-based features
- **File Size:** 126 MB (optimized)
- **Format:** CSV (UTF-8 encoding)
- **Quality Score:** 1.0/1.0 (perfect validation)

Technical Coverage

- **Frequency Range:** 300-600 GHz (complete THz spectrum)
- **Distance Range:** 1m to 5km (lab to wide-area)
- **SINR Range:** -30 to +50 dB (realistic conditions)
- **Throughput Range:** 0.1 to 1000 Gbps (6G targets)
- **Latency Range:** 0.011 to 2.841 ms (ultra-low latency)

Environmental Conditions

- **Clear Weather:** 127,122 samples (47.1%)
 - **Light Rain:** 57,129 samples (21.2%)
 - **Fog Conditions:** 40,795 samples (15.1%)
 - **Heavy Rain:** 30,876 samples (11.4%)
 - **Snow Conditions:** 14,078 samples (5.2%)
-

Parameter Categories

Atmospheric Turbulence (5 parameters)

- **Cn² Structure Parameter:** 1e-17 to 1e-13 m^{2/3}
- **Fried Parameter:** 0.01 to 1.0 m
- **Scintillation Index:** 0.001 to 2.0
- **Beam Wander:** ±0.1 to ±10 mrad

Weather & Environment (4 parameters)

- **Temperature:** -40°C to +60°C
- **Humidity:** 10% to 95%
- **Pressure:** 800 to 1200 hPa
- **Rain Rate:** 0 to 50 mm/h
- **Wind Speed:** 0 to 30 m/s

Fading & Channel (4 parameters)

- **Rician K-Factor:** 0 to 20 dB
- **Doppler Frequency:** 0 to 1000 Hz
- **Path Loss:** 60 to 180 dB
- **RMS Delay Spread:** 1 to 500 ns

Pointing & Alignment (4 parameters)

- **Beam Width:** 0.1 to 10 mrad
- **Pointing Error:** 0 to 5 mrad
- **Tracking Error:** 0 to 3 mrad
- **Pointing Loss:** 0 to 10 dB

OAM Beam Physics (4 parameters)

- **OAM Mode ():** -10 to +10
- **Mode Purity:** 0.7 to 0.99
- **Beam Divergence:** 0.5 to 5 mrad
- **Inter-Mode Crosstalk:** -40 to -10 dB

Hardware Impairments (3 parameters)

- **Phase Noise:** -120 to -80 dBc/Hz
- **I/Q Amplitude Imbalance:** 0 to 5 dB
- **Amplifier Efficiency:** 20% to 80%

Propagation Physics (4 parameters)

- **Wavelength:** 0.5 to 1.0 mm
- **Fresnel Radius:** 0.1 to 10 m
- **Diffraction Loss:** 0 to 20 dB
- **Oxygen Absorption:** 0.1 to 50 dB/km

Performance Metrics (5 parameters)

- **Throughput:** 0.1 to 1000 Gbps
- **Latency:** 0.011 to 2.841 ms
- **SINR:** -30 to +50 dB
- **Reliability:** 90.0% to 99.999%
- **Energy Efficiency:** 10x to 100x over 5G

File Structure

dataset/

generators/dataset/processed/

final-dataset-RL-processing.csv	# Main production dataset (270K)
comprehensive_dataset_270k.csv	# Full version with metadata
lab_controlled_comprehensive.csv	# Lab scenarios (70K)
indoor_realistic_comprehensive.csv	# Indoor scenarios (80K)
outdoor_urban_comprehensive.csv	# Urban scenarios (70K)
high_mobility_comprehensive.csv	# Mobility scenarios (50K)

metadata/

parameter_definitions.json	# Complete parameter docs
validation_report.json	# Quality metrics

generation_config.yaml	# Production config
publication_visuals/	
comprehensive_dataset_infographic.png	# Parameter overview
system_architecture_diagram.png	# System design
scenario_coverage_matrix.png	# Coverage matrix
dataset_workflow_flowchart.png	# Generation process
research_impact_diagram.png	# Research impact

Quick Start

Load Dataset

```
import pandas as pd
import numpy as np

# Load main dataset
df = pd.read_csv('generators/dataset/processed/comprehensive-correlated-6g-Dataset.csv')

print(f"Dataset shape: {df.shape}")
print(f"Scenarios: {df['scenario'].value_counts()}")
print(f"Parameter columns: {len(df.columns)} total")
```

Basic Analysis

```
# Key performance metrics
print("Performance Ranges:")
print(f"Throughput: {df['throughput_gbps'].min():.1f} - {df['throughput_gbps'].max():.1f} Gbps")
print(f"Latency: {df['latency_ms'].min():.3f} - {df['latency_ms'].max():.3f} ms")
print(f"SINR: {df['sinr_db'].min():.1f} - {df['sinr_db'].max():.1f} dB")
print(f"Reliability: {df['reliability_percent'].min():.1f}% - {df['reliability_percent'].max():.1f}%")

# Scenario distribution
scenario_stats = df.groupby('scenario').agg({
    'throughput_gbps': ['mean', 'std'],
    'latency_ms': ['mean', 'std'],
    'sinr_db': ['mean', 'std']
}).round(2)
print(scenario_stats)
```

Visualization

```
import matplotlib.pyplot as plt
import seaborn as sns

# Performance distribution by scenario
fig, axes = plt.subplots(2, 2, figsize=(15, 10))

# Throughput by scenario
```

```

sns.boxplot(data=df, x='scenario', y='throughput_gbps', ax=axes[0,0])
axes[0,0].set_title('Throughput Distribution by Scenario')
axes[0,0].tick_params(axis='x', rotation=45)

# Latency by scenario
sns.boxplot(data=df, x='scenario', y='latency_ms', ax=axes[0,1])
axes[0,1].set_title('Latency Distribution by Scenario')
axes[0,1].tick_params(axis='x', rotation=45)

# SINR by scenario
sns.boxplot(data=df, x='scenario', y='sinr_db', ax=axes[1,0])
axes[1,0].set_title('SINR Distribution by Scenario')
axes[1,0].tick_params(axis='x', rotation=45)

# Weather condition impact
weather_throughput = df.groupby('weather_condition')['throughput_gbps'].mean()
weather_throughput.plot(kind='bar', ax=axes[1,1], title='Throughput by Weather')
axes[1,1].tick_params(axis='x', rotation=45)

plt.tight_layout()
plt.show()

```

Validation Results

Quality Metrics

- **Overall Quality Score:** 1.0/1.0 (Perfect)
- **Parameter Coverage:** 100% of specified ranges
- **Physics Consistency:** Strong correlations (0.43+)
- **Statistical Validity:** Normal distributions where expected
- **No Invalid Values:** Zero out-of-range entries

Applications

Deep Reinforcement Learning

- Multi-agent network optimization
- Dynamic beam steering control
- Adaptive power allocation
- Interference mitigation strategies

6G System Design

- OAM mode selection optimization
- Channel prediction and compensation
- Network slicing parameter tuning
- Coverage and capacity planning

Research Applications

- THz propagation modeling validation
- Atmospheric effect quantification
- Hardware impairment characterization
- Standards compliance testing

Industry Applications

- Equipment testing and validation
 - Network deployment planning
 - Performance prediction models
 - Regulatory compliance assessment
-

Competition Instructions and Challenges

Overview

The 6G OAM THz Dataset is designed to enable comprehensive machine learning competitions and research challenges. This section provides detailed instructions for organizing competitions, evaluation metrics, and analysis tools to ensure fair and meaningful comparisons across different approaches.

Competition Framework

Track 1: Throughput Optimization Challenge **Objective:** Maximize data throughput while maintaining QoS constraints

Dataset Split: - **Training Set:** 200,000 samples (Lab: 50K, Indoor: 60K, Urban: 50K, Mobility: 40K) - **Validation Set:** 40,000 samples (Lab: 10K, Indoor: 10K, Urban: 10K, Mobility: 10K) - **Test Set:** 30,000 samples (Lab: 10K, Indoor: 10K, Urban: 10K) - **Hidden labels**

Evaluation Metrics: - **Primary:** Mean throughput improvement (%) - **Secondary:** Latency penalty coefficient - **Constraint:** SINR > -10 dB, Reliability > 95%

Baseline Algorithm:

```
# Simple linear regression baseline
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error, r2_score

def baseline_model(X_train, y_train, X_test):
    model = LinearRegression()
    model.fit(X_train, y_train)
    y_pred = model.predict(X_test)
    return y_pred

# Load dataset
df_train = pd.read_csv('training_split.csv')
features = ['sinr_db', 'distance_m', 'frequency_ghz', 'oam_mode', 'weather_condition']
```



```
X_train = df_train[features]
y_train = df_train['throughput_gbps']
```

Baseline performance: ~15.3 RMSE, $R^2 = 0.67$

Track 2: Multi-Objective Optimization Challenge **Objective:** Optimize throughput, latency, and energy efficiency simultaneously

Evaluation Function:

```
def competition_score(throughput, latency, energy_efficiency):
    # Normalize metrics to 0-1 scale
    norm_throughput = (throughput - throughput_min) / (throughput_max - throughput_min)
    norm_latency = 1 - (latency - latency_min) / (latency_max - latency_min) # Lower is better
    norm_energy = (energy_efficiency - energy_min) / (energy_max - energy_min)

    # Weighted combination
    score = 0.4 * norm_throughput + 0.3 * norm_latency + 0.3 * norm_energy
    return score
```

Pareto Frontier Analysis:

```
import matplotlib.pyplot as plt
from sklearn.metrics import mean_absolute_error

def plot_pareto_frontier(results_df):
    fig, ax = plt.subplots(1, 3, figsize=(15, 5))

    # Throughput vs Latency
    ax[0].scatter(results_df['throughput'], results_df['latency'], alpha=0.6)
    ax[0].set_xlabel('Throughput (Gbps)')
    ax[0].set_ylabel('Latency (ms)')
    ax[0].set_title('Throughput vs Latency Trade-off')

    # Throughput vs Energy
    ax[1].scatter(results_df['throughput'], results_df['energy_efficiency'], alpha=0.6)
    ax[1].set_xlabel('Throughput (Gbps)')
    ax[1].set_ylabel('Energy Efficiency (x over 5G)')
    ax[1].set_title('Throughput vs Energy Efficiency')

    # Latency vs Energy
    ax[2].scatter(results_df['latency'], results_df['energy_efficiency'], alpha=0.6)
    ax[2].set_xlabel('Latency (ms)')
    ax[2].set_ylabel('Energy Efficiency (x over 5G)')
    ax[2].set_title('Latency vs Energy Efficiency')

    plt.tight_layout()
    plt.show()
```

Track 3: Deep Reinforcement Learning Challenge **Objective:** Design DRL agents for dynamic network optimization

Environment Setup:

```
import gym
import numpy as np

class THz6GEnvironment(gym.Env):
    def __init__(self, dataset_path):
        super(THz6GEnvironment, self).__init__()
        self.df = pd.read_csv(dataset_path)

        # State space: 33 physics parameters
        self.observation_space = gym.spaces.Box(
            low=-np.inf, high=np.inf, shape=(33,), dtype=np.float32
        )

        # Action space: [beam_angle, power_level, oam_mode]
        self.action_space = gym.spaces.Box(
            low=np.array([-10, 0.1, -10]),
            high=np.array([10, 1.0, 10]),
            dtype=np.float32
        )

        self.current_idx = 0
        self.episode_length = 500

    def reset(self):
        self.current_idx = np.random.randint(0, len(self.df) - self.episode_length)
        return self._get_observation()

    def step(self, action):
        # Apply action and get reward
        reward = self._calculate_reward(action)
        self.current_idx += 1

        done = (self.current_idx % self.episode_length == 0)
        obs = self._get_observation()

        return obs, reward, done, {}

    def _calculate_reward(self, action):
        # Multi-objective reward: throughput + latency penalty + energy bonus
        current_state = self.df.iloc[self.current_idx]

        throughput_bonus = current_state['throughput_gbps'] / 1000 # Normalize
        latency_penalty = -current_state['latency_ms'] / 10 # Penalty
        energy_bonus = current_state['energy_efficiency'] / 100 # Bonus
```

```
    return throughput_bonus + latency_penalty + energy_bonus
```

Submission Requirements

Code Submission

```
# Required structure for submissions
class CompetitionSolution:
    def __init__(self):
        """Initialize your model/algorithm"""
        pass

    def fit(self, X_train, y_train):
        """Train your model on training data"""
        pass

    def predict(self, X_test):
        """Generate predictions for test data"""
        pass

    def get_model_info(self):
        """Return model description and hyperparameters"""
        return {
            'model_name': 'YourModelName',
            'hyperparameters': {},
            'training_time': 0.0,
            'memory_usage': 0.0
        }

# Example submission
solution = CompetitionSolution()
solution.fit(X_train, y_train)
predictions = solution.predict(X_test)
```

Evaluation Script

```
def evaluate_submission(submission_file, test_data_path):
    """
    Evaluate competition submission

    Args:
        submission_file: Path to submission CSV
        test_data_path: Path to hidden test set

    Returns:
        dict: Evaluation metrics
    """
    import pandas as pd
```

```

from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score

# Load predictions and ground truth
predictions = pd.read_csv(submission_file)
ground_truth = pd.read_csv(test_data_path)

# Calculate metrics
mse = mean_squared_error(ground_truth['target'], predictions['prediction'])
mae = mean_absolute_error(ground_truth['target'], predictions['prediction'])
r2 = r2_score(ground_truth['target'], predictions['prediction'])

# Competition-specific metrics
throughput_improvement = calculate_throughput_improvement(predictions, ground_truth)
constraint_violations = calculate_constraint_violations(predictions, ground_truth)

return {
    'mse': mse,
    'mae': mae,
    'r2_score': r2,
    'throughput_improvement': throughput_improvement,
    'constraint_violations': constraint_violations,
    'final_score': calculate_final_score(mse, throughput_improvement, constraint_violations)
}

```

Analysis Tools and Utilities

Dataset Analysis Toolkit

```

class DatasetAnalyzer:
    def __init__(self, dataset_path):
        self.df = pd.read_csv(dataset_path)

    def scenario_analysis(self):
        """Analyze performance across different scenarios"""
        scenario_stats = self.df.groupby('scenario').agg({
            'throughput_gbps': ['mean', 'std', 'min', 'max'],
            'latency_ms': ['mean', 'std', 'min', 'max'],
            'sinr_db': ['mean', 'std', 'min', 'max']
        }).round(3)
        return scenario_stats

    def correlation_analysis(self):
        """Analyze parameter correlations"""
        import seaborn as sns

        correlation_matrix = self.df.corr()
        plt.figure(figsize=(15, 12))
        sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', center=0)

```

```

plt.title('Parameter Correlation Matrix')
plt.show()

return correlation_matrix

def weather_impact_analysis(self):
    """Analyze weather condition impacts"""
    weather_analysis = self.df.groupby('weather_condition').agg({
        'throughput_gbps': 'mean',
        'latency_ms': 'mean',
        'path_loss_db': 'mean'
    }).round(3)

    # Visualization
    fig, axes = plt.subplots(1, 3, figsize=(18, 5))
    weather_analysis['throughput_gbps'].plot(kind='bar', ax=axes[0],
                                             title='Throughput by Weather')
    weather_analysis['latency_ms'].plot(kind='bar', ax=axes[1],
                                         title='Latency by Weather')
    weather_analysis['path_loss_db'].plot(kind='bar', ax=axes[2],
                                           title='Path Loss by Weather')

    plt.tight_layout()
    plt.show()

    return weather_analysis

```

Performance Benchmarking

```

def benchmark_algorithms():
    """Benchmark different algorithms on the dataset"""
    from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor
    from sklearn.neural_network import MLPRegressor
    from sklearn.svm import SVR

    algorithms = {
        'Linear Regression': LinearRegression(),
        'Random Forest': RandomForestRegressor(n_estimators=100, random_state=42),
        'Gradient Boosting': GradientBoostingRegressor(n_estimators=100, random_state=42),
        'Neural Network': MLPRegressor(hidden_layer_sizes=(100, 50), random_state=42),
        'SVM': SVR(kernel='rbf')
    }

    results = {}
    for name, algorithm in algorithms.items():
        start_time = time.time()
        algorithm.fit(X_train, y_train)
        training_time = time.time() - start_time

```

```

y_pred = algorithm.predict(X_test)
mse = mean_squared_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)

results[name] = {
    'mse': mse,
    'r2_score': r2,
    'training_time': training_time
}

return pd.DataFrame(results).T

```

Competition Timeline and Prizes

Timeline (Example)

- **Week 1-2:** Dataset release and baseline publication
- **Week 3-8:** Development phase with public leaderboard
- **Week 9:** Final submissions deadline
- **Week 10:** Evaluation and winner announcement

Prize Categories

1. **Overall Winner:** Best performance across all metrics
2. **Innovation Award:** Most novel approach or methodology
3. **Efficiency Award:** Best performance-to-computation ratio
4. **Reproducibility Award:** Best documented and reproducible solution

Submission Format

```

submission/
  code/
    main.py           # Main algorithm implementation
    requirements.txt  # Dependencies
    README.md         # Algorithm description
  results/
    predictions.csv   # Test set predictions
    validation_scores.json
  documentation/
    methodology.pdf   # Technical approach
    results_analysis.md # Performance analysis
  reproducibility/
    training_log.txt  # Training process log
    hyperparameters.json

```

Getting Started

Quick Competition Setup

License & Usage

This dataset is released under the **MIT License**, enabling both academic and commercial use with proper attribution.

Usage Rights

- Academic research and publication
- Commercial product development
- Standards development and testing
- Educational and training purposes

Attribution Requirements

- Cite the dataset in publications
 - Link to the original repository
 - Notify authors of significant use cases
-

Contact

Dataset Creators: - Principal Investigator: Srivatsa Davuluri - **Email:** connect.davuluri@gmail.com
- **GitHub:** [@srivatsadavuluriiii](#)

Repository: [6G-OAM-THz-Dataset](#)