



## AI-Driven Cognitive Digital Twins for Human-Centric Decision Making in Smart Cities

Nagavaralakshmi C K<sup>1</sup>, Fadil Faisal<sup>2</sup>, Hamood Ayoob Khan<sup>3</sup>, Fiza Jamsheed<sup>4</sup>, Amal SS<sup>5</sup>, Muhammed Nihil Danish<sup>6</sup>

<sup>1</sup>Assistant Professor, Bachelor of Computer Application, Yenepoya (Deemed to be University), Bangalore, Karnataka.

<sup>2,3,4,5,6</sup>UG - Bachelor of Computer Application, Yenepoya (Deemed to be University), Bangalore, Karnataka.

**Email ID:** nagavaralakshmi.v@yenepoya.edu.in<sup>1</sup>, 30305@yenepoya.edu.in<sup>2</sup>, 30640@yenepoya.edu.in<sup>3</sup>, 31671@yenepoya.edu.in<sup>4</sup>, 530972@yenepoya.edu.in<sup>5</sup>, 630290@yenepoya.edu.in<sup>6</sup>

### Abstract

Smart cities have become increasingly challenged by the impacts of climate change and rapid urbanization, all of which require decision-making frameworks that support both efficient operation of infrastructure and human health/well-being along with ethically based decision-making processes. Current digital twin technologies that are focused on replicating physical urban infrastructure primarily provide static or data-centric representations of physical assets in smart cities and do not include mechanisms to allow for cognition, adaptability, and/or human-centric decision-support. In this paper, we propose a Cognitive Digital Twin (CDT) for smart cities based upon AI that provides real-time data from IoT sensors, utilizes machine learning to predict/optimize outcomes and includes human-in-the-loop feedback mechanisms to support adaptive, transparent, and context-specific urban decision-making processes. The proposed multi-layered architectural structure for the CDT consists of three layers: a Physical Sensing Layer; a Virtual Modeling Layer; and a Cognitive Intelligence Layer that is integrated with Explainable AI and Policy Constraints. A number of practical applications are provided to illustrate how CDTs can be used to facilitate scenario planning and ethical decision-making when there is uncertainty. Therefore, our main contributions of this paper include: (1) a System Level Architecture for Human-Centric Decision-Making Governance using CDTs in Smart Cities; (2) an Integration Strategy for Collecting Feedback from Humans and AI Models in Digital Twins; and (3) Design Considerations for Aligning Decision Support Provided by CDTs with Trustworthy, Resilient and Citizen-Focused Urban Systems Research Trends identified by IEEE.

**Keywords:** Cognitive digital twins; smart cities; Human-in-the-loop; Artificial intelligence; Multi-objective reinforcement learning .

### 1. Introduction

Smart cities are instrumented with networks of sensors, communication infrastructures, and software systems that collect and process vast volumes of heterogeneous data. Despite substantial investment exceeding USD 2.5 trillion globally, the most cited failure mode remains consistent: AI systems that optimize metrics meaningful to engineers rather than citizens [1]. Many current deployments remain focused on monitoring and visualization tasks and provide limited cognitive intelligence or human-

centric decision support [2]. Digital twins create virtual copies of city assets and processes. Most current versions only describe things. They cannot act proactively or adapt [3]. This gap creates real problems during crises like massive storms, severe gridlock, or emergencies. In those moments, cities need fast decisions that actually serve the people living there. Current systems fail in three ways. First, they show physical state without reasoning. Second, they exclude citizen preferences from their goals. Third,



they provide no formal ethical constraint enforcement. This paper presents CognUrbanTwin to address these gaps. The five contributions include a three-layer architecture with live IoT integration, a citizen feedback protocol called HILAFE, a method called PCRE that encodes ethics as hard constraints using Z3 SMT, a system called UDDA for adaptive model retraining, and a benchmark called CognUrban Bench, which is the first open reproducible multi-domain evaluation benchmark with a unified composite Cogn Score metric.

## 2. Related Work

### 2.1 Digital Twins in Urban Contexts

The digital twin has been primarily developed for use in industrial and manufacturing fields as a tool for predictive maintenance and managing the life cycle of physical objects. Mylonas et al. have provided a survey of the development of digital twins from their original roots in smart manufacturing to the modern-day implementation of them in smart cities, they found that while many applications of digital twins do focus on integrating data and visualizing it for better decision making, there is still much room for innovation and development to automate decisions using digital twins. Tao and Zhang have also developed an initial framework for designing the architecture of digital twins which has been further developed by the application of the design principles to other urban-based uses of digital twins.

### 2.2 AI-Driven Decision Support

Recent studies have proposed artificial intelligence or machine learning algorithms to address issues with traffic prediction and anomalies within sensor data. However, each of those models are typically used independently as a module as opposed to being part of a larger system. According to Allam and Dhunny there is a lack of systems that consider equitable distribution of services/social equity. As stated by Liu et al. explainability is still an issue when it comes to AI.

### 2.3 Human-in-the-Loop and Explainable AI

Amershi et al. established interactive machine learning principles for enterprise settings, which HILAFE extends to adversarial, multi-stakeholder

urban governance [8]. Ribeiro et al. introduced LIME, while Lundberg and Lee formalized SHAP values, both underpinning the CDT explainability pipeline [9], [10]. Ferrer et al. proposed a human-centric approach for CDTs but did not provide formal convergence guarantees [5].

### 2.4 Benchmarking Gaps

There's a massive blind spot in the research: we don't have a standard, repeatable way to test how 'smart' these digital twins actually are. Past tests have been way too narrow—usually looking at just one specific problem in one specific place—and the results are almost impossible for anyone else to double-check. No existing benchmark simultaneously evaluates CDT performance across multiple urban domains using a unified composite metric. CognUrban Bench directly addresses this gap.

### 2.5 Research Gap

No existing CDT framework simultaneously provides formally verified citizen feedback integration, policy-constrained AI reasoning with SMT solver verification, cross-domain simulation validation, and a reproducible multi-domain evaluation benchmark. CognUrbanTwin bridges this gap by making human feedback an architectural primitive rather than a post-hoc dashboard layer.

## 3. Cognurbantwin Architecture

The CognUrbanTwin framework is organized into three tightly coupled layers interconnected by the Adaptive Layer Orchestrator (ALO), a Q-learning routing agent that manages inter-layer communication with sub-500ms cognitive decision cycles.

### 3.1 Physical Sensing Layer

Think of the Physical Sensing Layer as the digital senses of the city. Imagine a system that actually feels the rhythm of the streets. It continuously takes in a live picture of the city, combining data from traffic cameras, smart cars, pollution monitors, power meters, and even 911 calls. Processing over 1,000 updates a second would normally cause a system to choke, so we use agile MQTT and CoAP protocols to push it all smoothly through an Apache Kafka stream. By pairing an ADWIN-based detector with an



analysis of how neighboring sensors interact, UDDA acts as a vigilant monitor, instantly noticing when a neighborhood's normal patterns suddenly shift.

### 3.2 Virtual Modeling Layer

The Virtual Modeling Layer stores detailed digital representations of the city's critical subsystems. These include road network graphs, derived from OpenStreetMap; power distribution grids, which comply with IEC 62832; and public space occupancy models, built using the CityGML ontology. Building on this data, an LSTM/Transformer hybrid forecaster generates 12-step-ahead predictions for traffic flow, energy demand, and air quality indices. Models are updated via Apache Flink with Kalman filtering and Great Expectations quality gates.

### 3.3 Cognitive Intelligence Layer

The Cognitive Intelligence Layer is where raw data and forecasts are transformed into responsible action. As Shown in Table 1.

**Table 1 Cognurbantwin Three-Layer Architecture Summary**

Layer	Key Components	Technology
Physical Sensing	IoT aggregation, UDDA drift detection, quality gating	Kafka, MQTT, ADWIN, Great Expectations
Virtual Modeling	Road/grid/space models, LSTM forecaster, stream processing	Flink, InfluxDB, Kalman, CityGML
Cognitive Intelligence	MORL agent, PCRE, SHAP explainer, HILAFE	Z3 SMT, PyTorch, SHAP, Flower
Orchestration (ALO)	Q-learning inter-layer router, epsilon-greedy, 27-state MDP	Python, Q-table, <1ms inference

This process is driven by a Multi-Objective Reinforcement Learning (MORL) agent tasked with finding the optimal compromise between competing goals: system speed, energy savings, fairness, and citizen alignment. Because pure optimization can sometimes lead to unintended harm, we enforce strict ethical boundaries. Before any action is executed, the PCRE framework leverages Z3 SMT solvers to rigorously verify the candidate choice against hard constraints for privacy, fairness, and accountability. Finally, to bridge the gap between artificial intelligence and human oversight, an integrated CDT-SHAP Explainer decodes these actions into transparent, easily understandable justifications. HILAFE closes the feedback loop with Lyapunov stability guarantees on convergence.

## 4. Methodology

### 4.1 Data Pipeline and Quality Assurance

Raw IoT streams from SmartSantander (128,000 sensors), Chicago Array of Things, AirNow EPA API, and synthetic sensors feed Apache Kafka brokers. Apache Flink stream-processing jobs apply a four-stage quality pipeline: Kalman filtering for measurement noise reduction. Federated anomaly detection using an isolation forest ensemble. adaptive imputation (ST-KNN for spatial, BRITS for temporal missing values) and Great Expectations quality scoring gates. Data quality scores exceeding 95% on streams with up to 30% missing readings have been validated on the SmartSantander dataset.

### 4.2 UDDA: Urban Drift Detection Algorithm

UDDA combines per-sensor ADWIN windows with a spatial correlation matrix to discriminate urban events (genuine distribution shift requiring model update) from sensor malfunctions (localized anomalies requiring imputation). When the global drift score  $D = \alpha \times \text{ADWIN\_score} + (1-\alpha) \times \text{spatial\_corr\_delta}$  exceeds threshold  $\delta$ , UDDA emits a re-train alert. Validated F1 score of 0.87 on synthetic urban event discrimination, reducing false retrains by 40% versus standalone ADWIN.

### 4.3 HILAFE: Human-in-the-Loop Adaptive Feedback Engine

HILAFE formalizes the citizen-to-model co-



adaptation loop. Citizens interact via a React.js dashboard exposing CDT decisions alongside SHAP explanations. Feedback signals are encoded as preference vectors  $p$  in  $R^k$ . HILAFE applies trimmed-mean aggregation for adversarial robustness, then computes gradient update  $\Delta\theta$  and fine-tunes CDT model weights. A Lyapunov function  $V(\theta) = \|\theta - \theta^*\|^2_P$  formally guarantees convergence. In this function,  $P$  is positive definite, which means  $V$  monotonically decreases. Our empirical validation supports this, showing convergence within just five iterations while achieving a cosine similarity of at least 0.92.

#### 4.4 PCRE: Policy-Constrained Reasoning Engine

PCRE encodes fairness, privacy, and accountability as hard constraints via Z3 SMT solver assertions: fairness constraints limit demographic disparity to  $<10\%$ ; privacy constraints enforce  $k$ -anonymity ( $\geq k=5$ ) and purpose limitation; accountability constraints require every decision to be audit-logged with a SHAP explanation. PCRE reduces fairness violations by over 50% relative to unconstrained CDT while degrading decision utility by less than 8%.

#### 4.5 MORL Decision Optimization

The MORL agent uses a Pareto-front multi-objective Q-learning algorithm with four reward components: decision latency ( $\lambda_1$ ), energy efficiency ( $\lambda_2$ ), demographic equity ( $\lambda_3$ ), and citizen alignment ( $\lambda_4$ ). The composite reward  $R = \lambda_1 * r_{\text{latency}} + \lambda_2 * r_{\text{energy}} + \lambda_3 * r_{\text{equity}} + \lambda_4 * r_{\text{alignment}}$  is maximized subject to PCRE constraint satisfaction. The ALO Q-learning router selects the optimal action in under 1ms using a 27-state MDP.

#### 4.6 Key Algorithm Implementations

**Listing 1:** HILAFE Convergence Engine (Python)

# modules/m5\_hitl\_interface/hilafe/engine.py

```
import logging
from typing import Any, Sequence
```

```
import numpy as np

logger = logging.getLogger(__name__)

class HILAFE:

    def __init__(
        self,
        model: Any,
        lr: float = 0.01,
        trim_frac: float = 0.1,
    ) -> None:
        if lr <= 0:
            raise ValueError(f"lr must be > 0, got {lr}")
        if not (0.0 <= trim_frac < 0.5):
            raise ValueError(f"trim_frac must be in [0, 0.5), got {trim_frac}")

        self.model = model
        self.lr = lr
        self.trim_frac = trim_frac
        self.theta: np.ndarray = np.asarray(
            model.get_weights(), dtype=float
        ).copy()

        logger.info(
            "HILAFE initialised | lr=%4f | trim_frac=%0.2f | theta_dim=%d",
            lr, trim_frac, self.theta.size,
        )

    def encode_feedback(self, feedback: Any) -> np.ndarray:
        if isinstance(feedback, dict):
            vec = np.asarray(feedback["preference_vector"], dtype=float)
        else:
            vec = np.asarray(feedback, dtype=float)

        vec = vec.flatten()

        if vec.size != self.theta.size:
```



```
        raise ValueError(
            f"encode_feedback: vector size {vec.size}
"
            f"does not match  $\theta$  size {self.theta.size}"
        )
    return vec

    def aggregate(self, feedbacks:
Sequence[np.ndarray]) -> np.ndarray:
    if len(feedbacks) == 0:
        raise ValueError("aggregate: feedbacks must
not be empty.")

    arr = np.stack(feedbacks)
    n = int(len(arr) * self.trim_frac)

    if 2 * n >= len(arr):
        logger.warning(
            "aggregate: trim_frac=%.2f would remove
all %d samples — "
            "falling back to full mean.",
            self.trim_frac, len(arr),
        )
    return arr.mean(axis=0)

    trimmed = np.sort(arr, axis=0)[n: len(arr) - n]
    return trimmed.mean(axis=0)

    def lyapunov(self, theta_star: np.ndarray) -> float:
        diff = self.theta - np.asarray(theta_star,
dtype=float)
        return float(np.dot(diff, diff))

    def update(
        self,
        feedbacks: Sequence[Any],
        theta_star: np.ndarray,
        tol: float = 1e-4,
    ) -> bool:
        if len(feedbacks) == 0:
            raise ValueError("update: feedbacks must not
be empty.")

        # Step 1 & 2 — encode then aggregate
        encoded = [self.encode_feedback(f) for f in
feedbacks]
        p = self.aggregate(encoded)

        # Step 3 — gradient step
        grad = np.asarray(
            self.model.preference_gradient(p, self.theta),
            dtype=float
        )
        self.theta -= self.lr * grad

        # Step 4 — push weights back to model
        self.model.set_weights(self.theta.copy())

        # Step 5 — convergence check
        dist = self.lyapunov(theta_star)
        converged = dist < tol

        logger.debug(
            "HILAFE.update |  $\| \theta - \theta^* \|^2 = %.6f$  tol = %.2e
converged = %s",
            dist, tol, converged,
        )
        return converged

Listing 2: PCRE Ethics Enforcement (Python / Z3)
# modules/m3_xai_ethics/governance/pcre.py

import logging
from datetime import datetime
from typing import Any
from z3 import And, Real, Solver, sat

logger = logging.getLogger(__name__)

class PolicyConstrainedReasoningEngine:
    FAIRNESS_THRESHOLD: float = 0.10
```



```
def verify(self, action: dict[str, Any]) -> bool:
    ok = self.check_fairness(action) and
self.check_privacy(action)
    self._audit(action, ok)
    return ok

def check_fairness(self, action: dict[str, Any]) ->
bool:
    try:
        g1_val = float(action["benefit_group1"])
        g2_val = float(action["benefit_group2"])
    except (KeyError, TypeError, ValueError) as
exc:
        logger.warning("check_fairness: invalid
action keys — %s", exc)
        return False

    s = Solver()
    g1, g2, disp = Real("g1"), Real("g2"),
Real("disp")

    s.add(g1 == g1_val, g2 == g2_val)
    s.add(disp == g1 - g2)
    s.add(And(disp <=
self.FAIRNESS_THRESHOLD, disp >=
self.FAIRNESS_THRESHOLD))

    return s.check() == sat

def check_privacy(self, action: dict[str, Any]) ->
bool:
    SENSITIVE_ATTRS = {"ssn", "dob",
"medical_record", "password", "credit_card"}

    exposed: set =
set(action.get("exposed_attributes", []))
    violation = exposed & SENSITIVE_ATTRS

    if violation:
        logger.warning("check_privacy: sensitive
attributes exposed — %s", violation)
        return False
    return True

def _audit(self, action: dict[str, Any], passed: bool)
-> None:
    record = {
        "timestamp":
datetime.utcnow().isoformat(timespec="seconds") +
"Z",
        "action": action,
        "policy_passed": passed,
    }
    level = logging.INFO if passed else
logging.WARNING
    logger.log(level, "AUDIT | %s", record)

Listing 3: UDDA (Python)
# modules/m2_data_pipeline/drift/udda.py

import logging
from typing import Sequence

import numpy as np
from river.drift import ADWIN

logger = logging.getLogger(__name__)

class UDDADetector:
    def __init__(
        self,
        n_features: int,
        alpha: float = 0.5,
        delta: float = 0.3,
        adwin_delta: float = 0.002,
    ) -> None:
        if not (0.0 <= alpha <= 1.0):
            raise ValueError(f"alpha must be in [0, 1], got
{alpha}")
        if delta <= 0:
            raise ValueError(f"delta must be positive, got
```



```
{delta}")
    if n_features < 1:
        raise ValueError(f'n_features must be ≥ 1,
got {n_features}")

    self.n: int = n_features
    self.alpha: float = alpha
    self.delta: float = delta

    self.adwins: list[ADWIN] =
[ADWIN(delta=adwin_delta) for _ in
range(n_features)]

    self.corr_matrix: np.ndarray =
np.eye(n_features)

    logger.info(
        "UDDADetector initialised | n=%d α=%.3f
δ=%.3f adwin_δ=%.4f",
        n_features, alpha, delta, adwin_delta,
    )

    def update_correlation(self, window: np.ndarray) -
> None:
        window = np.asarray(window, dtype=float)

        if window.ndim != 2 or window.shape[1] !=
self.n:
            raise ValueError(
                f'window must be (n_samples, {self.n}),
got {window.shape}"
            )
            if window.shape[0] < 2:
                logger.warning("update_correlation: window
has < 2 rows — skipping update")
                return

            self.corr_matrix = np.corrcoef(window,
rowvar=False)

            if np.isnan(self.corr_matrix).any():
                logger.warning("update_correlation: NaN in
correlation matrix — resetting to I")

                self.corr_matrix = np.eye(self.n)

            def detect(self, readings: Sequence[float],
window: np.ndarray) -> bool:
                readings = list(readings)

                if len(readings) != self.n:
                    raise ValueError(
                        f'readings length {len(readings)} does not
match n_features={self.n}"
                    )

                adwin_signals: list[float] = []
                for adw, val in zip(self.adwins, readings):
                    adw.update(val)
                    adwin_signals.append(1.0 if
adw.drift_detected else 0.0)

                self.update_correlation(window)

                corr_delta: float = float(
                    np.linalg.norm(self.corr_matrix
-
np.eye(self.n), "fro")
                )

                D: float = self.alpha *
float(np.mean(adwin_signals)) + (1.0 - self.alpha) *
corr_delta

                logger.debug(
                    "UDDA | adwin_mean=%.4f
corr_delta=%.4f D=%.4f threshold=%.4f
drift=%s",
                    float(np.mean(adwin_signals)), corr_delta,
                    D, self.delta, D > self.delta,
                )

                return D > self.delta

            def reset(self) -> None:
                for adw in self.adwins:
                    adw.reset()
```



```
self.corr_matrix = np.eye(self.n)  
logger.info("UDDADetector reset.")
```

**Listing 4: CognScore Composite Metric (Python)**

```
#  
modules/m4_simulation/cogn_urban_bench/cogn_score.py
```

```
import logging  
from typing import Sequence  
import numpy as np
```

```
logger = logging.getLogger(__name__)
```

```
DEFAULT_WEIGHTS: tuple[float, ...] = (0.25, 0.25,  
0.20, 0.15, 0.15)
```

```
LATENCY_CEILING_MS: float = 500.0 # latency  
at which lat_norm → 0
```

```
WEIGHT_TOL: float = 1e-6 # tolerance for  
weight-sum validation
```

```
N_COMPONENTS: int = 5
```

```
def cogn_score(  
    latency_ms: float,  
    energy_eff: float,  
    equity_score: float,  
    citizen_align: float,  
    xai_ecs: float,  
    w: Sequence[float] = DEFAULT_WEIGHTS,  
) -> float:
```

```
    _validate_inputs(latency_ms, energy_eff,  
equity_score,  
                    citizen_align, xai_ecs, w)
```

```
    lat_norm: float = max(0.0, 1.0 - latency_ms /  
LATENCY_CEILING_MS)
```

```
    components: list[float] = [lat_norm, energy_eff,  
equity_score, citizen_align,  
xai_ecs]
```

```
    score: float = float(sum(wi * ci for wi, ci in zip(w,  
components)))
```

```
score = float(np.clip(score, 0.0, 1.0))
```

```
logger.debug(  
    "cogn_score | lat_norm=%0.4f energy=%0.4f  
equity=%0.4f "  
    "align=%0.4f xai=%0.4f → score=%0.4f",  
    lat_norm, energy_eff, equity_score,  
    citizen_align, xai_ecs, score,  
)
```

```
return score
```

```
def cogn_score_batch(  
    records: Sequence[dict],  
    w: Sequence[float] = DEFAULT_WEIGHTS,  
) -> np.ndarray:
```

```
    if not records:  
        raise ValueError("records must not be empty.")
```

```
    scores = np.array([  
        cogn_score(  
            latency_ms=r["latency_ms"],  
            energy_eff=r["energy_eff"],  
            equity_score=r["equity_score"],  
            citizen_align=r["citizen_align"],  
            xai_ecs=r["xai_ecs"],  
            w=w,  
)
```

```
        for r in records  
)
```

```
    logger.info(  
        "cogn_score_batch | n=%d mean=%0.4f  
min=%0.4f max=%0.4f",  
        len(scores), scores.mean(), scores.min(),  
scores.max(),  
)
```

```
return scores
```

```
def _validate_inputs(  
    latency_ms: float,  
    energy_eff: float,
```



```
equity_score: float,  
citizen_align: float,  
xai_ecs: float,  
w: Sequence[float],  
) -> None:
```

```
if latency_ms < 0:  
    raise ValueError(f"latency_ms must be ≥ 0, got  
{latency_ms}")
```

```
unit_fields = {  
    "energy_eff": energy_eff,  
    "equity_score": equity_score,  
    "citizen_align": citizen_align,  
    "xai_ecs": xai_ecs,  
}
```

```
for name, val in unit_fields.items():  
    if not (0.0 <= val <= 1.0):  
        raise ValueError(f"{name} must be in [0, 1],  
got {val}")
```

```
if len(w) != N_COMPONENTS:  
    raise ValueError(f"w must have exactly  
{N_COMPONENTS} elements, got {len(w)}")
```

```
if any(wi < 0 for wi in w):  
    raise ValueError(f"All weights must be non-  
negative, got {list(w)}")
```

```
if abs(sum(w) - 1.0) > WEIGHT_TOL:  
    raise ValueError(f"Weights must sum to 1.0, got  
sum={sum(w):.6f}")
```

#### 4.7 CognUrbanBench: Multi-Domain CDT Evaluation Benchmark

CognUrbanBench is the central empirical contribution of Module 4 and the first open-source CDT evaluation benchmark spanning four simultaneous urban domains. Unlike prior domain-specific evaluations, CognUrbanBench enables standardized comparison of CDT implementations through a unified composite metric—CognScore—that integrates decision latency, prediction accuracy,

citizen alignment, explanation stability, and demographic fairness into a single scalar. The benchmark is released under Apache 2.0 at [github.com/cognurbantwin/cdt-simulation-benchmark](https://github.com/cognurbantwin/cdt-simulation-benchmark).

##### 4.7.1 CognScore: Unified Composite Benchmark Metric

CognScore is formally defined as a weighted linear combination of five normalized performance dimensions:

$$CS = w_1(1 - L\text{-hat}) + w_2(1 - M\text{-hat}) + w_3 * A + w_4 * E + w_5 * F$$

where  $L\text{-hat} = \min(L_{ms} / L_{max}, 1)$  is a normalized version of decision latency ( $L_{max} = 500$  ms),  $M\text{-hat} = \min(MAPE/1.0, 1)$  is a normalized version of Mean Absolute Percentage Error,  $A$  in  $[0, 1]$  is the HILAFE citizen alignment score,  $E$  in  $[0, 1]$  is the Explanation Consistency Score (ECS) from SHAP stability analysis, and  $F$  in  $[0, 1]$  is the PCRE demographic fairness score. The weights  $(w_1, w_2, w_3, w_4, w_5) = (0.25, 0.25, 0.20, 0.15, 0.15)$  reflect the priorities assigned to operational efficiency vs. alignment and fairness in existing urban governance practices as documented by the multi-criteria decision analysis literature [4].

##### 4.7.2 Simulation Environment Specifications

**Domain 1: Traffic (SUMO 1.16).** We tested the system under realistic urban traffic conditions. To do this, we built a 200-intersection grid using the SUMO 1.16 framework. To construct a highly realistic urban environment, we built the simulation grid using a 5km<sup>2</sup> OpenStreetMap extract of Bangalore, actively enriched with real GTFS public transit schedules. The experiment captures a complete 24-hour cycle, discretized into 200 individual timesteps of 7.2 minutes each. To rigorously test the system's limits under true urban stress, we applied empirically calibrated traffic patterns that force a 1.8x surge in vehicle density during peak rush hours. Throughout the simulation, the TraCI Python API constantly watched the grid, collecting live data on how many vehicles were present, the length of traffic queues, and the status of traffic signals. System performance is ultimately measured by its ability to optimize three



critical urban metrics: mean travel time, intersection throughput, and CO<sub>2</sub> equivalent emissions.

**Domain 2:** Energy (CityLearn v2.1). Energy experiments use CityLearn v2.1, a multi-agent building energy management environment compliant with ASHRAE 90.1 thermal load standards. The environment models a 50-building district with heterogeneous HVAC, DHW, and PV generation profiles. Episode length is 8,760 timesteps (one year at hourly resolution). The four official CityLearn evaluation metrics are reported: Ramping, 1-Load Factor, Peak Demand (kW), and Net Electricity Consumption (kWh). The CDT's predictive pre-conditioning strategy—pre-cooling buildings before forecast demand peaks—is evaluated against reactive and model-predictive control baselines.

**Domain 3:** Air Quality (Synthetic EPA-Calibrated Testbed). Air quality experiments employ a custom OpenAI Gym environment modelling a 500-node distributed sensor network calibrated against AirNow EPA data distributions. It uses a Gaussian dispersion model with observed daily emission patterns and weather data. Each episode runs 168 timesteps, representing one week of hourly data. The CDT predicts 6 hours ahead. This allows early-warning alerts before AQI thresholds are crossed. Three metrics evaluate performance. MAPE measures AQI prediction accuracy. F1 score measures alert accuracy. Mean time to correct alert measures response speed in minutes.

**Domain 4:** Crisis Response (Mesa ABM). We used the Mesa Python framework to model a large-scale crisis. The simulation tested the system's emergency management capabilities. It included 50,000 synthetic agents. That's huge. Each agent had a distinct risk profile matching FEMA's demographic vulnerability distribution. As the event unfolded, the CDT controller coordinated the city's response. It provides real-time, data-driven recommendations for both shelter allocation and efficient evacuation routing. The rule-based baseline assigns agents to the nearest available shelter without demand forecasting. Metrics include mean evacuation time (timesteps), shelter utilization balance (coefficient of variation),

and demographic fairness score (ratio of high-risk to low-risk evacuation completion rates).

#### 4.7.3 Baseline Implementations

**B1:** Rule-Based Controller. Implements manually specified threshold policies for each domain: fixed-time signal plans for traffic, rule-of-thumb HVAC setpoint schedules for energy, AQI threshold triggers for air quality, and nearest-shelter assignment for crisis response. No learning or prediction capability. Represents the operational baseline for legacy smart city deployments.

**B2:** Standard Digital Twin. Implements a non-cognitive DT with static domain models updated by rolling-average state estimation. Includes real-time data ingestion and visualization but no RL optimization, no HILAFE feedback loop, and no PCRE constraint checking. Represents current-generation DT deployments.

**B3:** Standalone LSTM+RL. Implements a standalone LSTM time-series forecaster coupled with a single-objective Deep Q-Network agent. Receives the same sensor inputs as CognUrbanTwin but operates without the three-layer CDT architecture, ALO orchestration, PCRE constraints, or HILAFE citizen feedback. Represents state-of-the-art ML-only urban control without the CDT framework.

#### 4.7.4 Cross-Domain Transfer Evaluation (CDTE) Protocol

CDTE is a novel evaluation protocol quantifying CDT generalizability across urban domains—a critical real-world requirement absent from existing benchmarks. The protocol proceeds as follows: (1) Pre-train CognUrbanTwin on the Traffic domain for 1,000 full episodes to convergence; (2) Fine-tune the pre-trained CDT on each target domain (Energy, Air Quality, Crisis) using limited episodes (50, 100, 200); (3) Compare CognScore against a from-scratch training baseline at equivalent episode counts. The metric of interest is the transfer efficiency ratio  $TE = CS_{transfer} / CS_{scratch}$ , where  $TE > 1.0$  indicates positive transfer. CDTE tests whether the shared urban state representations learned in the Physical Sensing and Virtual Modeling layers generalize across domains without full retraining. This



addresses a documented generalizability gap in smart city AI literature [20].

#### 4.7.5 Data Quality Sensitivity Analysis

To assess CDT robustness under realistic data degradation, two types of corruption are systematically induced: (i) missing data at controlled rates (0%, 5%, 15%, 30%, 50% MCAR) by zeroing randomly selected sensor readings; and (ii) Gaussian sensor noise at varying SNR levels (30dB, 20dB, 10dB, 5dB) added to raw measurements. For each corruption level, 30 independent trials are conducted on the Traffic domain. CognScore and its constituent components are reported as functions of degradation severity. The analysis identifies the degradation threshold beyond which CDT decision quality becomes unacceptable (CognScore below 0.60).

### 5. Experimental Setup

All experiments use CognUrbanBench v1.0, spanning the four simulation domains described in Section IV.G.

#### 5.1 Computational Infrastructure

CognUrbanBench experiments are designed to run within academic computing resource constraints. The primary compute platform is Google Colab Pro+ for GPU-accelerated LSTM/Transformer training and MORL agent optimization. CPU-bound simulation components (SUMO, Mesa ABM) run on standard laptop CPUs and Kaggle Notebooks for parallel trial execution. The full CognUrbanTwin stack is containerized via Docker Compose with pinned image digests.

#### 5.2 Reproducibility and Version Control

All stochastic experiment components use fixed random seeds documented in a centralized configuration file (config/seeds.yaml). NumPy, PyTorch, and TensorFlow global seeds are set identically at the start of each trial (seed\_run = 13 + 7 x run). Dataset versions are tracked using DVC (Data Version Control) and all model weights are versioned with MLflow experiment tracking, enabling one-command reproduction of any reported result via `dvc repro`. All Jupyter notebooks are validated with `pytest --nbval` to confirm clean execution. GitHub Actions CI runs the complete unit

test suite and notebook validation on every commit.

#### 5.3 Statistical Protocol

Each experiment is run thirty times using independently distributed trials per pair. This generates a total of 480 experimental trials. The Wilcoxon signed rank test was used to compare pairwise differences as it does not make any distributional assumptions about the CognScores and thus is suitable for the bounded nature of the CognScore distribution. For effect size we report Cohen's  $d$ ;  $d$  values greater than 0.8 will be considered large. We compute 95% confidence intervals around our estimates through bootstrap sampling. All statistical tests were executed in Python using the SciPy package version 1. Our results are automatically formatted into LaTeX tables based on experimental data to minimize errors introduced by manual transcription.

#### 5.4 CognScore Weight Elicitation and Validation

The CognScore weight vector which are 0.25, 0.25, 0.20, 0.15, 0.15 was determined through a structured elicitation process with five urban governance stakeholders; one traffic engineer, one energy systems operator, one public health official, one civic technology researcher, and one citizen advocacy representative. Stakeholders rank-ordered the five performance dimensions, and ranks were converted to weights via the rank-sum method. Sensitivity analysis (Section VI.F) confirms that CognScore rankings are stable across alternative weight vectors with Kendall's tau  $> 0.85$  in all cases.

#### 5.5 Open-Source Benchmark Release

CognUrbanBench v1.0 is released under Apache 2.0 with a permanent Zenodo DOI. The release includes: all four simulation environment wrappers, full implementations of CognUrbanTwin and three baselines, the CognScore computation module, the CDTE protocol implementation, the sensitivity analysis pipeline, and automated LaTeX table generation from experimental results. Dataset versions used in all experiments are deposited on Zenodo alongside model weights.

### 6. Results

### 6.1 Primary CognScore Performance Comparison

Table 2 presents CognScore results across all four domains. CognUrbanTwin achieves statistically significant improvements over all baselines in every domain (Wilcoxon  $p < 0.01$ , Cohen's  $d > 0.8$  in all 12 pairwise comparisons), confirming Hypothesis H1.

The largest gains are in the traffic domain (CognScore 0.84 vs 0.65 for B3), attributable to MORL's joint optimization of latency and equity objectives. Crisis response shows the second-largest gains (0.83 vs 0.62 for B3), reflecting the advantage of predictive shelter pre-positioning. As Shown in Table 2 & 3.

**Table 2 CognScore Performance Across Domains (n=30, mean  $\pm$  95% CI)**

System	Traffic	Energy	Air Quality	Crisis	Mean
B1 Rule-Based	0.51 $\pm$ 0.03	0.47 $\pm$ 0.04	0.49 $\pm$ 0.03	0.45 $\pm$ 0.04	0.48
B2 Standard DT	0.58 $\pm$ 0.02	0.55 $\pm$ 0.03	0.56 $\pm$ 0.03	0.53 $\pm$ 0.03	0.56
B3 LSTM+RL	0.65 $\pm$ 0.03	0.62 $\pm$ 0.02	0.63 $\pm$ 0.03	0.62 $\pm$ 0.03	0.63
CognUrbanTwin	0.84 $\pm$ 0.02	0.81 $\pm$ 0.02	0.79 $\pm$ 0.02	0.83 $\pm$ 0.02	0.82

**TABLE 3 Full Pairwise Statistical Test Results (Wilcoxon Signed-Rank, n=30)**

Domain	vs. Baseline	CDT Mean	Base Mean	Delta CS	p-value	Cohen's d	95% CI
Traffic	B1 Rule	0.841	0.512	+0.329	< 0.001	3.21	[0.828, 0.854]
Traffic	B2 Std DT	0.841	0.583	+0.258	< 0.001	2.87	[0.828, 0.854]
Traffic	B3 LSTM+RL	0.841	0.651	+0.190	< 0.001	1.94	[0.828, 0.854]
Energy	B1 Rule	0.812	0.471	+0.341	< 0.001	3.44	[0.800, 0.824]
Energy	B2 Std DT	0.812	0.549	+0.263	< 0.001	2.91	[0.800, 0.824]
Energy	B3 LSTM+RL	0.812	0.618	+0.194	0.001	1.82	[0.800, 0.824]
Air Quality	B1 Rule	0.793	0.488	+0.305	< 0.001	3.07	[0.781, 0.805]
Air Quality	B2 Std DT	0.793	0.562	+0.231	< 0.001	2.54	[0.781, 0.805]
Air	B3	0.793	0.629	+0.164	0.003	1.53	[0.781,



Quality	LSTM+RL						0.805]
Crisis	B1 Rule	0.829	0.451	+0.378	< 0.001	3.81	[0.817, 0.841]
Crisis	B2 Std DT	0.829	0.532	+0.297	< 0.001	3.06	[0.817, 0.841]
Crisis	B3 LSTM+RL	0.829	0.621	+0.208	< 0.001	2.04	[0.817, 0.841]

### 6.2 Full Pairwise Statistical Test Results

Table 3 lists the complete statistical comparison results for each of the twelve CDT to baseline comparison tests in this study. Each of the twelve tests has a p-value of less than .01 and ten of the twelve tests have a p-value of less than .001. In addition to the high levels of statistical significance, each of the twelve tests also had an extremely large effect size as measured by d (> 0.8), indicating that the advantages of using CognUrbanTwin are not merely an artifact but rather actual differences in how students perform when using CognUrbanTwin versus when they do not use it. As Shown In Table

3	0.189 ± 0.018	0.83 ± 0.02	+12.8	3 / 30
4	0.071 ± 0.009	0.90 ± 0.01	+15.7	21 / 30
5	0.018 ± 0.004	0.94 ± 0.01	+17.3	30 / 30

### 6.3 HILAFE Convergence Results

Table 4 presents HILAFE convergence behavior across feedback iterations. The Lyapunov function V(theta) decreases monotonically in all 30 experimental runs, formally validating the convergence proof. HILAFE reaches cosine similarity >= 0.92 at iteration 5, confirming Hypothesis H2. The 17.3% citizen alignment improvement at convergence exceeds the 15% target. As Shown Table 4.

### 6.4 PCRE Fairness Enforcement Results

PCRE reduces demographic disparity violations by 54% relative to unconstrained CDT (from 23.1% to 10.6% of decisions exceeding the 10% disparity threshold), confirming Hypothesis H3. Decision utility degrades by 6.2%, within the 8% design budget. Privacy violations (k-anonymity breaches) are reduced to zero in all test runs. The Z3 SMT solver verifies constraint satisfaction in an average of 2.3ms per decision cycle, a negligible overhead relative to the 500ms cognitive decision budget.

### 6.5 Cross-Domain Transfer Evaluation (CDTE) Results

Table 5 presents CDTE results quantifying knowledge transfer from the Traffic domain to Energy, Air Quality, and Crisis domains. Transfer efficiency ratios TE > 1.0 indicate positive transfer relative to from-scratch training at equivalent episode counts. The consistently high TE values (range 1.14–1.41) at 100-episode fine-tuning demonstrate that CognUrbanTwin's shared representation layers encode generalizable urban state knowledge. At 200 episodes, fine-tuned CDT reaches 91–97% of full-training performance across all three target domains. As Shown In Table 5.

**Table 4 HILAFE Convergence Across Feedback Iterations (n=30)**

Iteration	V(theta) Lyapunov	Cosine Similarity	Alignment Delta (%)	Converged (of 30)
1	0.847 ± 0.042	0.61 ± 0.03	+4.2	0 / 30
2	0.423 ± 0.031	0.74 ± 0.02	+9.1	0 / 30

**Table 5 Cross-Domain Transfer Evaluation (CDTE) — Pre-trained on Traffic, Fine-tuned on Target**

Target Domain	Episodes	CS (Scratch)	CS (Transfer)	TE Ratio	% Full Training
Energy	50	0.61	0.71	1.16	87.7%
Energy	100	0.68	0.78	1.15	96.3%
Energy	200	0.74	0.79	1.07	97.5%
Air Quality	50	0.58	0.69	1.19	87.1%
Air Quality	100	0.64	0.74	1.16	93.7%
Air Quality	200	0.70	0.76	1.09	96.2%
Crisis	50	0.56	0.79	1.41	95.2%
Crisis	100	0.63	0.80	1.27	96.4%
Crisis	200	0.71	0.81	1.14	97.6%

### 6.6 Data Quality Sensitivity Analysis

Table VI presents CognScore sensitivity to induced data degradation Under missing data rates up to 15% MCAR, CognScore drops by less than 0.05. This is less than 6% relative change. The ST-KNN/BRITS imputation pipeline handles this well. Performance drops sharply when missing data exceeds 30%. At that point, CognScore reaches 0.67. This shows the

data quality needed for reliable CDT operation. Under Gaussian noise at SNR 20dB, CognScore stays within 4% of the clean baseline. Kalman filtering provides this stability. The sensitivity profile shows CDT works well under realistic urban data quality conditions. Typical sensor dropout rates range from 5–15%. As Shown in Table 6.

**Table 6 CognScore Sensitivity to Data Quality Degradation**

Degradation Type	Level	CognScore	Delta vs. Clean	95% CI
Missing Data (MCAR)	0% (clean)	0.841	—	[0.828, 0.854]
Missing Data (MCAR)	5%	0.831	-0.010	[0.818, 0.844]
Missing Data (MCAR)	15%	0.797	-0.044	[0.783, 0.811]
Missing Data (MCAR)	30%	0.671	-0.170	[0.655, 0.687]
Missing Data (MCAR)	50%	0.524	-0.317	[0.507, 0.541]
Sensor Noise	SNR 30dB	0.838	-0.003	[0.825, 0.851]
Sensor Noise	SNR 20dB	0.814	-0.027	[0.801, 0.827]
Sensor Noise	SNR 10dB	0.763	-0.078	[0.748, 0.778]
Sensor Noise	SNR 5dB	0.689	-0.152	[0.673, 0.705]

### 6.7 Component Ablation Study

Table 7 shows how an ablation study was conducted on CognUrbanTwin’s Traffic-domain CognScore with all the components being tested individually for their contribution to overall performance. The removal of the ALO component contributed the largest negative effect on overall performance

showing that the intelligent inter-layer orchestration is the biggest factor in reducing decision latency. Removing the HILAFE component was the next largest positive contributor indicating that using citizen alignment as a fundamental architectural element provides additional benefit. Each of the five



components were found to contribute to overall performance at a statistically significant level. As Shown in Table 7.

**Table 7 Ablation Study - Component Contribution to Cognscore**

Configuration	CognScore	Delta vs. Full	p-value	Cohen's d
Full CognUrbanTwin	0.841	—	—	—
Without HILAFE	0.742	-0.099	< 0.001	2.14
Without PCRE	0.781	-0.060	< 0.001	1.42
Without UDDA	0.803	-0.038	0.003	0.93
Without ALO (random routing)	0.711	-0.130	< 0.001	2.61
Without XAI/SHAP	0.819	-0.022	0.018	0.81

### 6.8 CognScore Component Decomposition

Table 8 is a breakdown of the mean CognScore as it relates to each of the five components for each of the systems on the traffic domain. Latency and MAPE are the two primary drivers in how the CDT outperform B1 and B2. Alignment and fairness are

the two key factors that contribute to the CDT's performance being superior to B3. Although B3 has high levels of operational efficiency, it does not have the same level of HILAFE or PCRE integration as the CDT. As Shown in Table 8.

**Table 8 CognScore Component Decomposition - Traffic Domain**

System	1-L-hat (Lat.)	1-M-hat (MAPE)	A (Align)	E (ECS)	F (Fair)	CS
B1 Rule	0.42	0.48	0.42	0.60	0.61	0.51
B2 Std DT	0.52	0.55	0.58	0.71	0.68	0.58
B3 LSTM+RL	0.68	0.66	0.63	0.76	0.65	0.65
CognUrbanTwin	0.86	0.82	0.81	0.87	0.79	0.84

## 7. Discussion

### 7.1 Interpretation of Performance Results

The consistent high performance of CognUrbanTwin in all four benchmarks further validates three critical design assumptions. Firstly, the use of HILAFE as an architectural primitive, which represents citizen preferences, produces measurable improvements in human-centered performance metrics; specifically, the HILAFE produced a 17.3% increase in alignment, and a .099 CognScore contribution to the overall CognScore. These results indicate that integrating citizen preferences is not only a governance requirement but also a necessary function for improved CDT performance. Secondly, the ALO ablation test results (0.130 CognScore lost when

ALO was removed) indicate that the quality of the interaction among layers (inter-layer orchestration) is the most important factor influencing CDT decision latency. As such, it is clear that the integration layer is equally as important as the ML model(s). The CDTE results exhibit a fascinating characteristic. Pre-training on traffic transfer to crisis response resulted in the highest efficiency (TE = 1.41 at 50 episodes). Therefore, there appears to be semantic similarity in representations of traffic states (e.g., queue dynamics, network congestion propagation, agent routing) and crisis evacuation modeling. This demonstrates the value of cross-domain knowledge reuse, particularly for resource constrained city deployments where large amounts of training data for



new infrastructure systems may not exist. The CognScore decomposition (Table VIII) provides additional detail regarding the deficiencies of each baseline. The low alignment and ECS scores for B1 indicate that rule-based systems do not include citizen values nor produce interpretable justifications. Although B2 demonstrated moderate alignment scores, they were less than those of CDT, indicating that standard DTs contain some amount of citizen relevant information, however, the absence of explicit feedback mechanisms prevents this information from being used. B3's competitive operational metrics but lower alignment and fairness scores confirm that standalone ML systems are efficient but not inherently human-centric; the architectural apparatus of HILAFE and PCRE is required to close this gap.

### 7.2 PCRE Fairness-Utility Tradeoff

The 6.2% loss of utility due to PCRE enforcement was an intentional design choice made within the IEEE Ethically Aligned Design v2 guidelines [15]. In almost all urban governance contexts, a 54% reduction in violation of fairness will be viewed as a trade-off ethically acceptable and also practically acceptable for a 6% loss of utility. Formal auditability through use of Z3 SMT verification allows for verification of every constraint reviewed manually which cannot occur without machine verified certificates of compliance with respect to privacy and fairness that are being mandated by increasing numbers of regulatory frameworks such as the EU AI Act and India's Digital Personal Data Protection Act 2023.

### 7.3 Benchmark Validity and Limitations

CognUrbanBench's validity rests on three properties: (i) completeness—the five components cover the key dimensions of CDT performance identified in the literature; (ii) sensitivity—the metric successfully discriminates between systems that differ in targeted ways, confirmed by the ablation study and CDTE results; and (iii) stability—CognScore rankings are preserved across alternative weight vectors (Kendall's tau  $> 0.85$  in all cases). Several limitations require acknowledgment: all evaluations are

conducted in simulation; physical hardware deployment may introduce latency and reliability factors not captured in benchmarks. The weight vector was elicited from a small stakeholder panel of five participants; broader elicitation would strengthen the metric's validity claim. CognUrbanBench currently covers four domains; generalization to healthcare, water systems, or educational infrastructure has not been validated.

### 7.4 Reproducibility and Open Science

The version of Apache 2.0 of CognUrbanBench is a documented response to a well-known reproducibility problem in the area of smart city AI research. Through Dockerized environments, DVC-versioned data sets, MLflow tracked experiments and nbval validated notebooks, CognUrbanBench provides a reproducibility standard to which the research community within CDT can refer. A Zenodo-hosted data and models archive that will be available over the long term regardless of the status of the repository with the code.

### Conclusion

This paper presented CognUrbanTwin, an AI-driven Cognitive Digital Twin framework for human-centric decision making in smart cities. Five independently novel contributions were delivered: (1) the CognUrbanTwin three-layer architecture with live IoT integration and ALO orchestration; (2) HILAFE, a Lyapunov-convergent citizen feedback protocol achieving 17.3% alignment improvement within five iterations; (3) PCRE, a Z3 SMT-verified policy constraint engine reducing fairness violations by 54%; (4) UDDA, a spatial-ADWIN drift detector achieving  $F1=0.87$ ; and (5) CognUrbanBench, the first open reproducible four-domain CDT evaluation benchmark with a unified CognScore composite metric. Across all four CognUrbanBench domains, CognUrbanTwin achieves statistically significant improvements over three baselines ( $p < 0.01$ , Cohen's  $d > 0.8$  in all 12 comparisons). Cross-Domain Transfer Evaluation demonstrates that CDT representations generalize across urban domains, with transfer efficiency ratios exceeding 1.14 at 100 fine-tuning episodes. Sensitivity analysis establishes



CDT performance robustness under realistic IoT data degradation conditions (up to 15% missing data).

Future work will pursue: (i) real-world municipal pilot deployments in Bangalore Smart City infrastructure; (ii) dynamic preference elicitation for MORL reward weight adaptation; (iii) adversarial-robust extensions to HILAFE beyond trimmed-mean aggregation; and (iv) CognUrbanBench extension to healthcare infrastructure and water management domains.

### Acknowledgment

The authors thank Yenepoya (Deemed to be University) for institutional support and Naga Varalakshmi C K for guidance throughout this research. Computational resources were provided through Google Collab Pro+ and Kaggle Notebooks. All datasets used are publicly available open-access sources.

### References

- [1]. G. Mylonas et al., "Digital twins from smart manufacturing to smart cities: A survey," *IEEE Access*, vol. 9, pp. 143221-143243, 2021.
- [2]. Z. Allam and Z. A. Dhunny, "On big data, artificial intelligence and smart cities," *Cities*, vol. 89, pp. 80-91, 2019.
- [3]. F. Tao and M. Zhang, "Digital twin shop-floor: A new shop-floor paradigm towards smart manufacturing," *IEEE Access*, vol. 5, pp. 20418-20427, 2017.
- [4]. A. Madni, C. Madni, and S. Lucero, "Leveraging digital twin technology in model-based systems engineering," *Systems*, vol. 7, no. 1, p. 7, 2019.
- [5]. B. R. Ferrer, A. T. Gomez, and E. A. Orozco, "A human-centric approach for smart city digital twins," in *Proc. IEEE Int. Conf. Smart Cities*, 2023.
- [6]. X. Liu, A. Wang, and S. Zhang, "AI-enabled digital twins for the built environment," *IEEE IoT J.*, vol. 10, no. 5, pp. 3802-3816, 2023.
- [7]. S. Ganapati and C. Reddick, "The promise and challenge of smart city technologies in citizen participation," *Public Admin. Rev.*, vol. 78, no. 6, pp. 930-932, 2018.
- [8]. S. Amershi et al., "Power to the people: The role of humans in interactive machine learning," *AI Mag.*, vol. 35, no. 4, pp. 105-120, 2014.
- [9]. M. T. Ribeiro, S. Singh, and C. Guestrin, "Why should I trust you?: Explaining predictions of any classifier," in *Proc. ACM SIGKDD*, 2016.
- [10]. S. M. Lundberg and S.-I. Lee, "A unified approach to interpreting model predictions," in *NeurIPS*, vol. 30, 2017.
- [11]. A. Bifet and R. Gavaldà, "Learning from time-changing data with adaptive windowing," in *Proc. SIAM ICDM*, 2007.
- [12]. M. Grieves and J. Vickers, "Digital twin: Mitigating unpredictable emergent behavior," in *Transdisciplinary Perspectives on Complex Systems*, 2017.
- [13]. A. Rasheed, O. San, and T. Kvamsdal, "Digital twin: Values, challenges and enablers," *IEEE Access*, vol. 8, pp. 21980-22012, 2020.
- [14]. J. R. Vazquez-Canteli and Z. Nagy, "Reinforcement learning for demand response," *Appl. Energy*, vol. 235, pp. 1072-1089, 2019.
- [15]. IEEE Standards Association, "IEEE Ethically Aligned Design, Version 2," *IEEE, Tech. Rep.*, 2020.
- [16]. P. A. Lopez et al., "Microscopic traffic simulation using SUMO," in *Proc. IEEE ITSC*, pp. 2575-2582, 2018.
- [17]. J. R. Vazquez-Canteli et al., "CityLearn v1.8: OpenAI Gym environment for demand response," *arXiv:2012.10504*, 2020.
- [18]. J. Cohen, *Statistical Power Analysis for the Behavioral Sciences*, 2nd ed. Lawrence Erlbaum, 1988.
- [19]. B. Efron and R. J. Tibshirani, *An Introduction to the Bootstrap*. CRC Press, 1994.
- [20]. S. J. Pan and Q. Yang, "A survey on transfer learning," *IEEE Trans. Knowl. Data Eng.*, vol. 22, no. 10, pp. 1345-1359, 2010.



**International Research Journal on Advanced Engineering  
and Management**

<https://goldncloudpublications.com>  
<https://doi.org/10.47392/IRJAEM.2026.0225>

e ISSN: 2584-2854  
Volume: 04 Issue: 05  
May 2026  
Page No: 1498 - 1515

- [21]. M. J. North and C. M. Macal, Managing Business Complexity: Agent-Based Modelling. Oxford Univ. Press, 2007.