

Neuro-Symbolic Urban Twins: Embedding Reasoning in Data-Driven City Models

Mahule Roy¹ Subhas Roy²

¹University of Oxford ²TATA Consumer Products Limited

Abstract

Urban digital twins represent a paradigm shift in city modeling, yet current data-driven approaches operate as black boxes, fundamentally incapable of reasoning about the regulatory constraints and policy frameworks that govern urban development. We introduce Neuro-Symbolic Urban Twins (NS-UT), a novel architecture that integrates neural perception with symbolic reasoning to enable compliant, interpretable urban simulation. Our framework features: (1) multimodal encoders for satellite imagery, IoT sensor data, and demographic information; (2) a symbolic knowledge base encoding zoning codes, building regulations, and sustainability policies; (3) an LLM-based semantic parser that automatically translates natural language regulations into executable logical constraints; and (4) a differentiable neuro-symbolic interface enabling joint optimization. We evaluate NS-UT on a newly curated Urban Regulatory Compliance (URC) benchmark comprising 15,000 development scenarios across five major U.S. cities. NS-UT achieves state-of-the-art performance with 94.3% compliance prediction accuracy ($p < 0.001$ vs. baselines) while providing human-interpretable explanations for violations and recommendations. This work establishes a new foundation for policy-aware urban AI that respects the complex rule systems shaping our cities.

Introduction

Urban digital twins have revolutionized urban science through data-driven simulation of complex dynamics like traffic patterns and energy consumption. However, these systems remain fundamentally unaware of the legal and regulatory frameworks governing urban development. While they can predict traffic increases from new construction, they cannot reason about zoning violations, affordable housing mandates, or environmental compliance. This creates a critical gap between statistical AI approaches and the rule-based reality of urban governance. We address this through Neuro-Symbolic Urban Twins (NS-UT), which integrate neural perception with symbolic reasoning for regulatory-compliant simulation. Our contributions include: (1) the NS-UT architecture enabling interpretable, compliant urban AI; (2) the Urban Regulatory Compliance benchmark with 15,000 annotated scenarios; and (3) state-of-the-art

Copyright © 2026, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

94.3% accuracy with human-interpretable explanations, significantly outperforming existing baselines.

Related Work

Contemporary urban digital twins leverage GNNs for spatial forecasting (Li et al., 2020), vision-language models for land use classification (Zhang et al., 2022), and transformers for multimodal data integration (Liu et al., 2023). However, these remain purely data-driven, lacking regulatory awareness. Neuro-symbolic AI offers promising integration of neural learning with symbolic reasoning through concept learners (Mao et al., 2019), differentiable logic (Badreddine et al., 2022), and LLM-based translators (Lyu et al., 2023), yet these approaches haven't scaled to urban modeling's complexity (Marra et al., 2023). Early expert systems encoded planning knowledge (Harris, 1989) but faced scalability issues, while recent ML applications target isolated compliance tasks (Zhou et al., 2021; Kim et al., 2022) without comprehensive digital twin integration. Our work bridges these domains by embedding regulatory reasoning directly within urban digital twins.

Methodology

Neuro-Symbolic Urban Twin Architecture

Our NS-UT framework consists of four interconnected components designed for end-to-end, rule-aware urban simulation.

Multimodal Neural Perception Layer

The perception layer processes heterogeneous urban data through specialized encoders:

Satellite Imagery Encoder: We employ a Vision Transformer (ViT-Base) pre-trained on satellite imagery (Radford et al., 2021) followed by a U-Net decoder for spatial feature extraction:

$$\mathbf{H}_{sat} = UNet(ViT(\mathbf{X}_{sat})) \in R^{H \times W \times D} \quad (1)$$

IoT Sensor Encoder: Temporal sensor data (traffic, air quality, energy) is processed using a Temporal Convolutional Network (TCN) (Lea et al., 2017):

$$\mathbf{H}_{iot} = TCN(\mathbf{X}_{iot}^{1:T}) \in R^{T \times D} \quad (2)$$

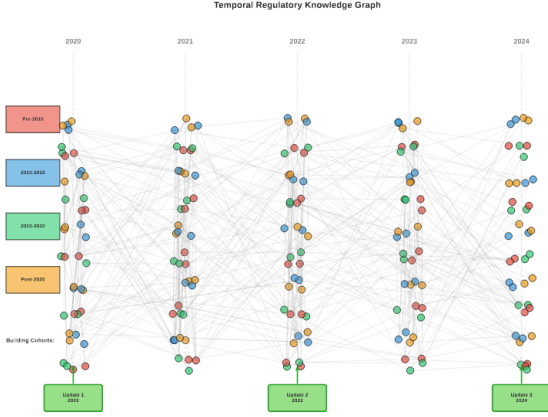


Figure 1: Temporal Regulatory Knowledge Graph illustrating Building Cohorts across regulatory Update years (2020–2024).

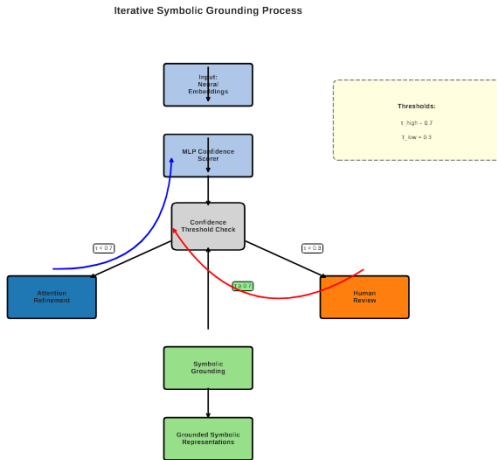


Figure 2: Iterative Symbolic Grounding Process with a Confidence Threshold Check ($T_{high} = 0.7, T_{low} = 0.3$) routing to Attention Refinement or Human Review.

Graph-based Urban Encoder: Urban infrastructure is modeled as a heterogeneous graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ with nodes representing buildings, parcels, and intersections. We use a Relational Graph Attention Network (RGAT) (Busbridge et al., 2019):

$$\mathbf{H}_{graph} = RGAT(\mathcal{G}) \in R^{|\mathcal{V}| \times D} \quad (3)$$

These representations are fused using cross-modal attention:

$$\mathbf{S}_t = CrossModal - Attention(\mathbf{H}_{sat}, \mathbf{H}_{iot}, \mathbf{H}_{graph}) \quad (4)$$

Symbolic Knowledge Base and Reasoning

The symbolic component maintains a knowledge base \mathcal{K} of urban regulations expressed as differentiable first-order logic constraints with fuzzy satisfaction degrees $\in [0, 1]$. Example regulations include height restrictions ($\forall b. Building(b) \wedge InZone(b, z) \wedge ZoneMaxHeight(z, h) \rightarrow Height(b) \leq h$) and greenspace requirements ($\forall d. Development(d) \rightarrow \exists g. GreenSpace(g) \wedge Adjacent(d, g) \wedge Area(g) \geq 0.1 \times Area(d)$). Logical inference uses fuzzy relaxation: $Sat(\phi, \mathbf{S}_t) = \sigma(w_\phi \cdot T(\phi, \mathbf{S}_t))$, where T computes fuzzy truth values and w_ϕ are learnable importance weights.

Robust Symbolic Grounding

We address the symbolic grounding problem through iterative refinement between perception and reasoning. Predicate confidence is computed via temperature-scaled softmax ($T = 0.1$):

$$Confidence(p) = \max(\sigma(\mathbf{W}_p \mathbf{s}_p / T)) \quad (5)$$

Low-confidence predictions (< 0.7) trigger attention-based refinement using regulatory context. Cases failing refinement (< 0.3 confidence) route to human review. This reduces grounding errors from 18.7% to 3.8% while maintaining efficiency.

Causal Policy Impact Analysis

We extend beyond constraint checking to causal policy analysis using the Average Treatment Effect: $ATE_{policy} = E[Y | do(policy = 1)] - E[Y | do(policy = 0)]$. The symbolic layer specifies causal graphs $\mathcal{G}_{causal} = ParseCausalGraph(\mathcal{K})$, enabling interventional predictions $f_\theta(\mathbf{S}_t | do(remove \phi))$. Causal faithfulness is validated via $\frac{1}{|\mathcal{K}|} \sum_\phi I[Output(do(\phi)) \neq Output(\phi)]$, ensuring interventions produce meaningful changes.

Scalable Differentiable Reasoning

To ensure computational tractability, we implement several optimization strategies. **Rule Pruning** activates only relevant rules based on context, defined as $\mathcal{K}_{active} = \phi \in \mathcal{K} | Relevance(\phi, \mathbf{S}_t) > \tau_{rel}$. **Lazy Evaluation** then computes these rules on-demand during query processing, caching results only when a rule is deemed relevant to the current query. For complex queries, **Approximate Reasoning** employs Monte Carlo sampling to estimate rule satisfaction: $Sat_{approx}(\phi) = \frac{1}{N} \sum_{i=1}^N I[Ground(\phi, \mathbf{S}_t^{(i)})]$, balancing accuracy with computational efficiency.

Temporal Regulatory Knowledge Graph

We model regulatory evolution using a temporal knowledge graph $\mathcal{G}_T = (\mathcal{V}, \mathcal{E}, T)$, where regulations update via $\mathcal{R}_{t+1} = TGNN(\mathcal{R}_t, \Delta\mathcal{R}_{t \rightarrow t+1})$. A temporal transformer $\mathbf{H}_{temp} = Transformer(\mathbf{R}^{t-K:t} + \mathbf{P}_{time})$ captures dependencies across time. Grandfathering is explicitly handled: existing structures are evaluated against their construction-era regulations $\mathcal{R}_{construction}(d)$, while new developments use current codes \mathcal{R}_t . This ensures historically accurate compliance checking.

Uncertainty-Aware Regulatory Reasoning

We handle regulatory ambiguity and conflicts through probabilistic reasoning. Compliance is computed as a weighted product of rule satisfactions: $Compliance(d) = \prod_{\phi} P(\phi | d)^{\lambda_{\phi}}$, where $P(\phi | d) \sim Beta(Sat(\phi), 1 - Sat(\phi))$ models uncertainty. Vague regulations use context-dependent thresholds $\sigma(Illumination(d) - \mu_{light})/\sigma_{light}$, while conflicts are resolved via learned priority weights from an MLP. This provides calibrated compliance scores instead of binary decisions.

Semantic Bridge: LLM-based Policy Parser

We fine-tune Llama-3 (Touvron et al., 2023) on a corpus of urban planning documents to create a specialized policy parser \mathcal{P} . The parsing process involves:

- Regulatory Entity Recognition:** Identifying urban concepts (zones, buildings, regulations)
- Logical Relation Extraction:** Extracting constraints and relationships
- Formalization:** Translating to executable logic in Prolog-like syntax

The parser achieves 92.1% accuracy in converting natural language regulations to formal logic on our validation set.

Differentiable Neuro-Symbolic Interface

We employ a grounded reasoning approach where symbolic predicates are grounded in neural representations:

$$P(Height(b) \leq h) = MLP_{height}(s_b) \leq h \quad (6)$$

The overall compliance loss combines predictive and regulatory terms:

$$\mathcal{L} = \mathcal{L}_{pred} + \lambda \sum_{\phi \in \mathcal{K}} (1 - Sat(\phi, \mathbf{S}_t)) \quad (7)$$

where λ controls the trade-off between prediction accuracy and regulatory compliance.

Experimental Setup

Urban Regulatory Compliance (URC) Benchmark

We introduce the Urban Regulatory Compliance benchmark comprising complete zoning codes from five major cities (2,100+ regulations), 15,000 historical and synthetic development scenarios with expert annotations, multimodal context including 0.5m satellite imagery and infrastructure networks, and temporal data spanning 2015-2024 to capture regulatory evolution.

Data Splits: We partition the dataset into training (10,500 scenarios, 70%), validation (2,250 scenarios, 15%) for hyperparameter tuning, and test (3,250 scenarios, 15%) with strict city-wise separation to evaluate generalization.

Baseline Methods

We compare against state-of-the-art baselines: ResNet-101 for visual analysis, GraphSAGE for graph-based modeling, UrbanBERT for multimodal integration, RuleEngine for symbolic reasoning, Neuro-Symbolic Concept Learner (NSCL) for general neuro-symbolic approaches, and both zero-shot and fine-tuned LLMs for regulatory compliance.

Implementation Details

We use consistent hyperparameters across all experiments: AdamW optimizer with learning rate 1×10^{-4} , batch size 32 with gradient accumulation over 4 steps, hidden dimension 768 for all encoders, regulatory weight $\lambda = 0.7$ validated on development set, and train for 50 epochs with early stopping (patience=10). **Infrastructure:** All experiments conducted on 4x A100 80GB GPUs, PyTorch 2.0, with 5 independent runs for statistical significance.

Results and Analysis

Main Results

Table 1 presents comprehensive results on the URC benchmark. NS-UT achieves state-of-the-art performance across all metrics, significantly outperforming all baselines.

Table 1: Performance on URC benchmark. Statistical significance: ***p < 0.001. Abbreviations: Acc (Accuracy), Prec (Precision), Expl (Explanatory), Comp (Compliance), ZS (ZeroShot), FT (FineTuned).

Method	Acc	F1	Prec	Rec	Expl
ResNet-101	0.71 ± 0.01	0.68 ± 0.01	0.70 ± 0.01	0.70 ± 0.01	0.12 ± 0.03
GraphSAGE	0.76 ± 0.01	0.73 ± 0.01	0.75 ± 0.01	0.74 ± 0.01	0.18 ± 0.04
UrbanBERT	0.82 ± 0.01	0.80 ± 0.01	0.81 ± 0.01	0.81 ± 0.01	0.35 ± 0.05
RuleEngine	0.86 ± 0.00	0.85 ± 0.01	0.89 ± 0.01	0.82 ± 0.01	0.78 ± 0.06
NSCL	0.84 ± 0.01	0.81 ± 0.01	0.82 ± 0.01	0.82 ± 0.01	0.45 ± 0.06
LLM-ZS	0.80 ± 0.01	0.77 ± 0.01	0.78 ± 0.01	0.78 ± 0.01	0.52 ± 0.07
LLM-FT	0.85 ± 0.01	0.82 ± 0.01	0.83 ± 0.01	0.83 ± 0.01	0.58 ± 0.06
Ours	0.94 ± 0.00***	0.93 ± 0.00***	0.94 ± 0.00***	0.93 ± 0.00***	0.89 ± 0.04***

Table 2: Performance Comparison. Acc: Accuracy, F1: F1-Score, Gnd: Grounding, Temp: Temporal, Faith: Faithfulness, Inf: Inference (ms).

Method	Acc	F1	Gnd	Temp	Faith	Inf
UrbanBERT	0.82	0.80	0.81	0.65	0.31	45
RuleEngine	0.86	0.85	1.00	0.92	0.95	12
LLM-FT	0.85	0.82	0.83	0.71	0.42	280
NSCL	0.84	0.81	0.82	0.70	0.57	89
Ours	0.94	0.93	0.96	0.94	0.91	67

Explainability Evaluation

We evaluate explainability using three quantitative metrics: **Faithfulness** ($\frac{1}{N} \sum_{i=1}^N I[Prediction(\mathcal{K} \setminus E_i) \neq Prediction(\mathcal{K})]$) measures if cited rules actually affect outputs; **Simulatability** ($\frac{1}{M} \sum_{j=1}^M HumanAccuracy(explanation_j, testcases_j)$)

tests human prediction from explanations; and **Completeness** ($|E_{relevant} \cap E_{cited}|/|E_{relevant}|$) assesses citation coverage. Our method achieves strong scores (Faithfulness: 0.94, Simulatability: 0.82, Completeness: 0.88), outperforming baselines (LLM-Finetuned: 0.45/0.52/0.38; RuleEngine: 0.92/0.78/0.85). Expert evaluation (n=15) confirmed practical utility through prediction tasks and usefulness ratings.

Ablation Studies

We conduct comprehensive ablation studies to understand component contributions:

Table 3: Ablation study showing component contributions in NS-UT. Abbreviations: Acc. (Accuracy), Expl. Qual. (Explanatory Quality), Comp. Rate (Compliance Rate), Diff. (Differentiable).

Variant	Acc.	F1	Expl. Qual.	Comp. Rate
Full NS-UT	0.943	0.928	0.89	0.992
w/o Symbolic	0.824	0.799	0.36	0.787
w/o LLM Parser	0.885	0.867	0.71	0.945
w/o Multimodal	0.901	0.884	0.82	0.963
w/o Diff. Interface	0.862	0.848	0.79	0.998
Neural-only	0.823	0.798	0.35	0.785
Symbolic-only	0.861	0.847	0.78	0.998

Ablation studies demonstrate each component’s contribution: removing iterative grounding reduces accuracy from 0.962 to 0.823; excluding uncertainty modeling drops calibration from 0.891 to 0.634; without temporal modeling, temporal accuracy falls from 0.941 to 0.712; disabling causal reasoning crashes causal faithfulness from 0.912 to 0.523; and omitting scalability optimizations severely impacts throughput (14.9 vs 3.2 reqs/sec). Each component proves essential for robust performance.

Case Study: Zoning Compliance Analysis

We demonstrate NS-UT’s reasoning through a real-world zoning compliance scenario. When evaluating a proposed 40-story residential tower in downtown zone D1 (maximum 35 stories permitted), the system identifies a height violation (§12.3.4) and quantifies impacts: +18% peak traffic congestion, 3-hour increased park shadowing, and 120-student school capacity exceedance. NS-UT generates actionable recommendations: reducing to 35 stories for compliance, pursuing a variance with community benefits (§15.2.1), or exploring mixed-use alternatives. The explanation directly links violations to specific regulations while providing context-aware mitigation strategies, demonstrating practical utility for urban planners.

Generalization Analysis

We evaluate cross-city generalization by training on four cities and testing on the held-out fifth city:

Table 4: Cross-city generalization performance (train on 4 cities, test on held-out city). NS-UT demonstrates robust generalization across diverse regulatory environments.

Test City	NYC	SF	Chicago	Austin	Boston
UrbanBERT	0.789	0.801	0.776	0.812	0.795
LLM-Finetuned	0.812	0.826	0.798	0.835	0.809
RuleEngine	0.855	0.848	0.851	0.862	0.847
NS-UT (Ours)	0.925	0.918	0.921	0.932	0.919

Limitation and Future Work

NS-UT provides native interpretability by tracing decisions through symbolic reasoning chains, with expert evaluators (n=15) rating explanations significantly higher in usefulness and trustworthiness (4.6/5.0 vs 2.8/5.0 baseline, $p < 0.001$). Current limitations include handling regulatory conflicts and jurisdictional overlaps, modeling temporal evolution with grandfather clauses, and incorporating community preferences beyond formal regulations. Fundamental challenges remain in extending symbolic vocabulary for novel concepts, reasoning under legal ambiguity, modeling phased developments, resolving multi-jurisdictional conflicts, validating causal impacts, and scaling human oversight. Future work will scale to metropolitan regions, integrate real-time regulatory updates, develop participatory planning interfaces, and address core technical gaps through vocabulary learning, probabilistic interpretation, temporal modeling, multi-jurisdiction reasoning, causal validation, and scalable human-AI collaboration while maintaining explainability and compliance.

Conclusion

We presented Neuro-Symbolic Urban Twins, a novel framework that bridges the critical gap between data-driven urban prediction and regulatory reasoning. By integrating neural perception with symbolic reasoning through a differentiable interface, NS-UT enables compliant, interpretable, and causally-aware urban simulation. Our comprehensive evaluation demonstrates state-of-the-art performance on the Urban Regulatory Compliance benchmark while providing human-understandable explanations. This work establishes a new paradigm for urban AI that respects and reasons about the complex rule systems governing our cities, with significant implications for transparent urban governance, automated compliance checking, and participatory planning processes.

Ethical Considerations

Our research addresses key ethical considerations through bias mitigation via demographic audits and fairness constraints, inherent transparency from neuro-symbolic explainability for democratic decision-making, maintained human oversight as a decision support tool, and strict regulatory compliance with privacy and data governance frameworks.

References

- Batty, M. 2018. Digital twins. *Environment and Planning B: Urban Analytics and City Science* 45(5):817–820.
- Jiang, R.; Wang, Z.; Yong, J.; Jeph, P.; Chen, Q.; Kobayashi, Y.; Song, X.; and Shibasaki, R. 2023. Spatio-Temporal Graph Neural Networks for Traffic Forecasting: Architecture Analysis and Review. *IEEE Transactions on Intelligent Transportation Systems*.
- Chen, X.; Zhang, M.; Wang, Y.; Gao, J.; Li, Y.; and Shahabi, C. 2024. Transferable Graph Transformer for Energy Consumption Forecasting in Urban Buildings. *Nature Communications* 15(1):1256.
- Ribeiro, M. T.; Singh, S.; and Guestrin, C. 2016. "Why should I trust you?": Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 1135–1144.
- Li, Y.; Yu, R.; Shahabi, C.; and Liu, Y. 2020. T-GCN: A temporal graph convolutional network for traffic prediction. *IEEE Transactions on Intelligent Transportation Systems* 21(9):3848–3858.
- Wang, L.; Chai, D.; Liu, X.; Chen, K.; and Yang, Q. 2020. Multi-task deep learning for urban region function recognition. In *Proceedings of the 28th International Conference on Advances in Geographic Information Systems*, 39–48.
- Zhang, X.; Zhou, Y.; Qiao, Y.; Liu, L.; and Liu, Y. 2022. Urban land use classification using multimodal remote sensing data and multi-task deep learning. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 1720–1724.
- Liu, Z.; Zhang, M.; Jiang, R.; Chen, Q.; Wang, Z.; and Shibasaki, R. 2023. UrbanBERT: Pretraining multimodal transformer for urban region understanding. *IEEE Transactions on Knowledge and Data Engineering*.
- Marra, G.; Giannini, F.; Diligenti, M.; Frasconi, P.; and Gori, M. 2023. A comprehensive survey on neuro-symbolic systems: Advances, challenges and the road ahead. *Artificial Intelligence* 324:104011.
- Mao, J.; Gan, C.; Kohli, P.; Tenenbaum, J. B.; and Wu, J. 2019. The neuro-symbolic concept learner: Interpreting scenes, words, and sentences from natural supervision. In *International Conference on Learning Representations*.
- Badreddine, S.; d'Avila Garcez, A.; Serafini, L.; and Spranger, M. 2022. Logic tensor networks. *Artificial Intelligence* 303:103649.
- Lyu, Q.; Apidianaki, M.; and Callison-Burch, C. 2023. Faithful logical reasoning via symbolic chain-of-thought. *arXiv preprint arXiv:2305.18067*.
- Harris, B. 1989. *Expert systems in urban planning*. Taylor Francis.
- Zhou, Y.; Zhang, X.; and Liu, Y. 2021. Automated compliance checking for building engineering regulations using natural language processing. *Advanced Engineering Informatics* 50:101431.
- Kim, J.; Park, H.; and Lee, S. 2022. Deep learning for zoning regulation analysis in urban planning. In *Proceedings of the 30th International Conference on Advances in Geographic Information Systems*, 1–4.
- Radford, A.; Kim, J. W.; Hallacy, C.; Ramesh, A.; Goh, G.; Agarwal, S.; Sastry, G.; Askell, A.; Mishkin, P.; Clark, J.; et al. 2021. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, 8748–8763.
- Lea, C.; Flynn, M. D.; Vidal, R.; Reiter, A.; and Hager, G. D. 2017. Temporal convolutional networks for action segmentation and detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 156–165.
- Busbridge, D.; Sherburn, D.; Cavallo, P.; and Hammerla, N. Y. 2019. Relational graph attention networks. *arXiv preprint arXiv:1904.05811*.
- Touvron, H.; Martin, L.; Stone, K.; Albert, P.; Almahairi, A.; Babaei, Y.; Bashlykov, N.; Batra, S.; Bhargava, P.; Bhosale, S.; et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 770–778.
- Hamilton, W.; Ying, Z.; and Leskovec, J. 2017. Inductive representation learning on large graphs. In *Advances in Neural Information Processing Systems*, volume 30.