GEO-INVARIANT LEAD SCORING WITH DOMAIN-ADVERSARIAL TRANSFORMERS

Anonymous authors

000

001

002003004

010 011

012

013

014

015

016

017

018

019

021

025

026

027 028 029

031

033

034

036

037

040

041

042

043

044

046

047

048

050 051

052

Paper under double-blind review

ABSTRACT

Predicting B2B lead conversion requires not only modeling long-range dependencies in richly sequenced customer interactions but also ensuring fair performance across under-represented geographies. While our DeepScore transformer backbone improved overall AUPR from 0.266 to 0.360, it exhibited significant geo-skew: majority-region (America) signals dominated feature learning (AUPR 0.474), leaving East-Asia (0.262) under-served. To address this, we embed a Domain-Adversarial Neural Network (DANN) module into DeepScore's architecture. A gradient-reversal layer connects multiple domain discriminators to the shared transformer encoder, enforcing a minimax game that drives hidden representations to be predictive of conversion outcomes while remaining uninformative across multiple demographic and firmographic domains. Simultaneously, lightweight geo-specific classifier heads capture region-specific conversion patterns while maintaining shared feature representations, preventing model divergence across geographic markets. DeepScore + geo-DANN achieves a 4.3% relative gain in macro-AUPR and reduces inter-region AUPR gaps by up to 12.3\%, all without degrading America accuracy. To our knowledge, this is the first demonstration of adversarial domain adaptation in large-scale B2B lead scoring, offering a scalable path to equitable, high-fidelity predictions across heterogeneous markets.

1 Introduction

Predicting whether a prospective B2B lead will convert, known as *lead scoring*, is essential for allocating scarce sales resources effectively. For instance, enterprise cloud services sales typically require significant investment of Account Executive time for multiple lengthy meetings and technical demos with senior decision-makers, while Solution Architects need to develop custom proofs-of-concept and address specific technical requirements. Modern pipelines generate long sequences of time-stamped interactions (emails, ad clicks, webinars, calls), which tree-ensemble methods (XGBoost, LightGBM) struggle to model without extensive feature engineering. Transformer-based architectures, by applying self-attention over all touchpoints, have recently matched or surpassed these baselines on structured and sequential tasks (Lim et al., 2021b; Author, 2025), simplifying preprocessing while capturing long-range dependencies.

However, a single global model trained on pooled data often *privileges majority regions* (e.g., North America) and under-serves lower-volume markets (e.g., East-Asia), yielding significant geo-skew in performance metrics. To address this, we introduce **Geo-DANN DeepScore**: a transformer backbone augmented with (i) a gradient-reversal layer connecting multiple domain discriminators that learn representations invariant to demographic and firmographic attributes, and (ii) lightweight per-region classifier heads that learn residual local patterns. On a dataset of 1.4M leads across 10 geographic markets, Geo-DANN DeepScore achieves a 4.3% relative gain in macro-AUPR and reduces interregion performance gaps by up to 12.3%, all without harming majority-region performance. This work is the first to apply adversarial domain adaptation at scale for B2B lead scoring, offering a practical blueprint for fair, high-fidelity predictions across heterogeneous markets.

1.1 MOTIVATION AND PRIOR APPROACHES

Traditional lead scoring systems have evolved from rule-based qualification to machine learning models. These systems at large multinational technology companies must handle millions of active

leads across diverse geographic regions and business units, with significant downstream impact on sales prioritization and resource allocation. A wrong decision may have a negative impact on the course of millions of dollars. The high stakes of these decisions, where leads are assigned priority grades (from high-priority requiring immediate sales engagement to lower-priority for automated nurturing) that directly influence sales engagement, demand both accuracy and fairness across all served markets.

Traditional lead scoring approaches have predominantly relied on region-specific model deployment strategies, wherein distinct predictive models are developed and maintained for individual geographic regions or business units. These conventional methods typically employ gradient boosting algorithms (such as LightGBM) and necessitate extensive feature engineering processes tailored to regional characteristics. While such region-specific modeling can effectively capture local market dynamics and behavioral patterns, this approach presents significant scalability and operational challenges. The maintenance of multiple parallel models substantially increases computational overhead, complicates feature pipeline architecture, and impedes systematic performance evaluation across regions. Furthermore, this methodology exhibits inherent limitations in data-sparse environments, where low-volume regions possess insufficient training samples to develop statistically robust standalone models. Consequently, emerging markets and smaller geographic segments often experience degraded predictive performance, creating inconsistencies in lead scoring quality across the organization's global footprint. These operational inefficiencies and performance disparities highlight the need for a more unified and adaptable modeling framework.

1.2 The DeepScore Architecture

To address these limitations, we developed **DeepScore**, a unified transformer-based architecture with three key design goals: (1) consolidate multiple regional models into a single, globally-deployable model, (2) leverage fine-grained sequential interaction data rather than aggregated tabular features, and (3) minimize reliance on manual feature engineering.

DeepScore processes customer interactions as sequences through a transformer encoder, analyzing both temporal patterns and textual content associated with each touchpoint. This approach captures the full customer journey, from initial website visits to email exchanges, without destroying sequential information through aggregation. Initial evaluations demonstrated substantial improvements over gradient boosting baselines (see Section 5.2).

However, analysis revealed significant regional performance disparities: high-volume regions with abundant training data saw improvements exceeding 60%, while lower-volume markets experienced gains below 20%. Such disparities stem from global loss minimization that naturally overfits to majority domains, a problem that standard class-imbalance techniques like SMOTE can actually exacerbate (Piccininni et al., 2024; van den Goorbergh et al., 2022; Carriero et al., 2025).

1.3 ALTERNATIVE APPROACHES AND THEIR LIMITATIONS

Multi-Head Architecture. We initially explored a multi-head variant inspired by multi-task learning, where a shared transformer encoder feeds \$G\$ independent region-specific classifiers. The hypothesis was that each head would capture regional patterns while the shared backbone would benefit from pooled data. However, this approach yielded mixed results with a net decrease in macro-AUPR.

Three factors explain this failure: (i) Gradient interference: Back-propagation from high-volume regions dominates updates, causing the shared encoder to preferentially encode majority-region statistics (Yu et al., 2020; Standley et al., 2020). (ii) Representation entanglement: Without explicit regularization against geographic information, the encoder embeds regional cues that prevent effective cross-region transfer (Bousmalis et al., 2016; Ganin et al., 2016b). (iii) Sample inefficiency: Region-specific heads receive limited mini-batches, leading to slow convergence and high gradient variance (Zamir et al., 2018; Ruder, 2017).

Given these limitations, we explored simpler rebalancing strategies before adopting the adversarial approach.

Loss Re-weighting Strategies. We also experimented with two re-weighting schemes: (i) inverse frequency weighting based on regional training data availability ($w_g = N_{total}/N_g$), and (ii) inverse performance weighting where regions are weighted proportionally to $1/\text{AUPR}_g$ from initial validation. As shown in Table 5.2, both approaches degraded overall performance while failing to close regional gaps.

The fundamental issue is that re-weighting amplifies gradients from underrepresented regions without addressing the core problem: the encoder still learns geography-specific patterns rather than invariant representations. We're essentially amplifying noisy signals from sparse data without learning features that actually transfer across domains. Re-weighting changes the optimization emphasis but doesn't induce the feature invariance necessary for true domain adaptation (Zhao et al., 2019a).

Class Rebalancing. Traditional rebalancing techniques (SMOTE, random over/undersampling) are poorly suited to our geographic imbalance. Our regional sample sizes vary significantly (Americas: 40% of examples, East-Asia: 15%). Downsampling to the minority region would discard over 60% of our training data—an unacceptable loss of signal. Upsampling minorities through duplication or synthetic generation (SMOTE) risks overfitting to limited behavioral patterns and can actually amplify biases present in small samples (van den Goorbergh et al., 2022). Moreover, geographic regions exhibit genuine distributional differences in business practices and sales cycles that synthetic examples cannot capture. These approaches fundamentally mistake the problem as one of sample size rather than distribution shift.

Other Simple Approaches. We considered several other standard techniques: (i) Focal loss (Lin et al., 2017) which down-weights easy examples but doesn't address geographic shift; (ii) Progressive fine-tuning where we train on majority regions then fine-tune on minorities, risking catastrophic forgetting; and (iii) Ensemble methods combining region-specific models, which reintroduces the operational complexity we sought to avoid. Each fails to achieve the key goal: learning a single model with invariant features that generalizes across all regions while maintaining operational simplicity.

1.4 Domain-Adversarial Solution

To address these shortcomings, we integrated Domain-Adversarial Neural Networks (DANN) into DeepScore. DANN employs a gradient-reversal layer coupled with a domain discriminator to learn features that are simultaneously predictive of conversion and invariant to geographic origin. This approach, originally developed for computer vision (Ganin et al., 2016b), has shown promise when combined with transformers in object detection (Zhang et al., 2024) and e-commerce applications (Herold et al., 2025).

Our contribution demonstrates that adversarial domain adaptation can effectively address geographic performance disparities in B2B lead scoring at scale, providing a principled solution to the fairness-accuracy trade-off in multi-region deployment scenarios.

Beyond predictive performance, DANN uniquely enables business intelligence: the discriminator monitors regional divergence while the feature decomposition identifies which patterns should be global best practices versus legitimate local adaptations, which are capabilities that re-weighting approaches cannot provide.

The key insight is that geographic performance gaps stem from *distribution shift*, not mere *class imbalance*. Re-weighting and resampling optimize the wrong objective, they change which examples matter most but don't change what features are learned. DANN directly optimizes for invariance through its minimax objective, forcing the encoder to discard geographic artifacts while preserving conversion-predictive signal. This principled approach to distribution alignment, rather than ad-hoc rebalancing, explains the consistent improvements we observe in Section 5.2.

2 RELATED WORK

2.1 THE EVOLUTION OF LEAD SCORING

Lead scoring has evolved from rule-based systems encoding sales heuristics (D'Haen et al., 2013) to modern machine learning approaches. The current industry standard commonly employs gradient

boosting methods, XGBoost (Chen & Guestrin, 2016) and LightGBM (Ke et al., 2017), due to their strong performance on tabular data. However, these methods require extensive feature engineering to capture temporal patterns. Engineers create hundreds of aggregated features across multiple time windows, fundamentally destroying the sequential nature of customer journeys (Kumar & Shah, 2018).

Recent work has explored deep learning for lead scoring, but primarily through simple neural networks on the same aggregated features (Zhang & Wang, 2019). The critical insight, that the journey itself is the signal, has been overlooked in favor of incrementally improving feature engineering.

2.2 Geo-Aware Modeling in Other Domains

Global platforms and scientific applications have adopted a variety of strategies to balance shared learning with local adaptation. In traffic forecasting, Spatio-Temporal Graph Convolutional Networks (STGCN) learn a single graph-based encoder over an entire road network, while more recent evolutionary GNNs dynamically update region-specific adjacency structures to capture local traffic patterns without training separate models for each city (Yu et al., 2018; Chen et al., 2024b). In large-scale recommendation systems, two-tower embedding architectures train a global user-item model that is then fine-tuned or re-indexed at the city level, reducing the need for thousands of city-specific models and lowering operational overhead (Wang et al., 2022)

2.3 Cross domain modeling

Environmental and spatial forecasting have employed multi-task frameworks with a shared backbone plus lightweight domain-specific heads, demonstrating improved spatial generalization while avoiding fully separate regional models (Liu et al., 2025). CDTrans combines self-attention pseudo labels to close domain gaps across diverse image datasets, offering a blueprint for single-model multi-domain deployment without per-domain retraining (Xu et al., 2022).

2.4 Domain-Adversarial Learning

Ganin et al. (2016b) introduced DANN, embedding gradient reversal within back-propagation so that feature extractors become domain-confusing while label predictors remain discriminative. Ben-David et al. (2010) formalized generalization bounds using the $\mathcal{H}\Delta\mathcal{H}$ divergence between source and target distributions. Later work extended DANN to vision transformers, object detection and remote-sensing segmentation, proving scalability to high-capacity backbones. However, Zhao et al. (2019b) showed that representation invariance alone can be insufficient without an optimal joint hypothesis, motivating our geo-specific heads.

2.5 Transformers for Heterogeneous Sequential Data

Extending transformers beyond text requires handling heterogeneous inputs and irregular timing. Temporal Fusion Transformers (Lim et al., 2021a) introduced specialized components for time series forecasting with static covariates. TabTransformer (Huang et al., 2020b) applies attention to tabular features but doesn't model true sequences. Most relevant is the line of work on multi-modal transformers (Tsai et al., 2019), though these typically handle aligned modalities (e.g., video and audio) rather than the diverse interaction types in customer journeys.

2.6 Transformers on Classification Tasks

Since their introduction for sequence modeling, transformers have been broadly adopted for classification across multiple modalities. In natural language processing, pretrained encoders such as BERT fine-tuned on GLUE benchmarks have set new state-of-the-art on sentence and document classification tasks Devlin et al. (2019a); Liu et al. (2019).

2.7 Transformers on Structured Data

TabTransformer first demonstrated that self-attention contextualizes categorical embeddings to rival tree-ensemble models on tabular tasks Huang et al. (2020a). Subsequent surveys chronicle dozens

of transformer-based variants for tabular representation learning Badaro et al. (2023); Somvanshi et al. (2024); Ruan et al. (2024); Gorishniy et al. (2021); Arík & Pfister (2021); Singh et al. (2023); Wang & Sun (2022); Chen et al. (2024a); Fan & Waldmann (2024), but none address domain shift in business-scale B2B marketing sequences at the scale tackled here.

220 221

222

3 METHODOLOGY

224 225

3.1 DEEPS

226227228229230

238239240241242243244

247248249250

245246

251252253254255256257

260261262263264265

266

267

268

269

258

259

3.1 DEEPSCORE BACKBONE

Interaction Features

Profile (Lead Features) Interaction-level Time-related features features (year, Relative Head Lead Lead (channel, type month, day, Position feature : feature 2 feature N etc.) weekday) Prediction Prediction for Geo 1 Conversion Path Transformer Encoder Classifier Domain GRLAdversarial Path Prediction Head Lead Interaction History Prediction for Geo N **Conversion Path** GRI Adversarial Path

Figure 1: DeepScore Architecture (main)

DeepScore turns the full marketing and sales history of a lead into one long sequence and feeds it to a Transformer encoder. Every interaction (*e.g.* e-mail open, ad click, webinar attendance, phone call) is first mapped to a *learned touch embedding*. This is done through learning feature level embeddings for the various metadata attributes for a touch, and encoding the textual context into a semantic embedding. These feature embeddings, semantic embeddings, and numerical values are concatenated and fed into a linear reduction layer, producing the *learned touch embedding*. Alongside those touch tokens, we concatenate four *discrete time embeddings* (year, month, day-of-month, weekday) so the model can reason about seasonality and working-day effects; borrowing from Lim et al. (2021b). We also adapt the bucketed bias introduced for T5 in to encode relative position because relative distances are more descriptive with with very large sequences Shaw et al. (2018). The result is a sequence of lead interactions in the following shape. [Touch, Time, RelPos]₁ [Touch, Time, RelPos]₂

We then learn a representation for static "profile" information about the lead, using 241 categorical and 37 numerical features. This is handled in a similar way to the touch embeddings, where we learn feature-level embeddings for the categorical features, concatenate them with numerical features and reduce them into a profile embedding. We concatenate the lead interaction sequence with the profile information using a separator token and then front pad each sequence to meet a fixed distance T. This information is fed into a transformer encoder, following Vaswani et al. (2017).

```
\underbrace{\text{pad}}_{\times T-k}, \underbrace{\text{touch/time/pos}}_{\times k}, [SEP], \underbrace{\text{profile}}_{1}]
```

The encoder produces hidden states hidden $\in \mathbb{R}^{T \times 256}$. As in BERT(Devlin et al., 2019b), the penultimate position is reserved as a [CLS]-style token that summarizes the sequence; the model selects that vector and treats it as the dense lead representation $f_{\theta}(x)$. A single linear layer then converts $f_{\theta}(x)$ into the logit of conversion, yielding the backbone's binary prediction.

3.2 Multi-head prediction design

To improve the models ability to focus on distinct geos, we replace the single linear classifier by a torch.nn.ModuleDict that stores ten single-layer MLP heads. During a forward pass we look up the geography ID for every lead, select the corresponding head, and compute its logits, where C is the individualized classification heads:

$$\text{logits}_i = C_{\phi_{g(i)}} \big(f_{\theta}(x_i) \big), \quad g(i) \in \{1, \dots, 10\}.$$

3.3 GEO-DANN - DOMAIN-ADVERSARIAL NEURAL NETWORK

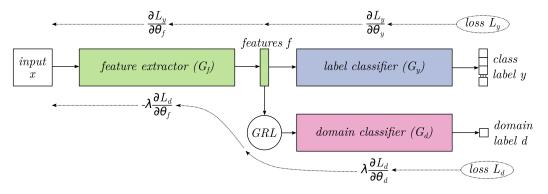


Figure 2: Information flow with a Gradient-Reversal Layer (GRL). During back-propagation the domain-loss gradient is sign-flipped, forcing the encoder to learn geography-invariant features while the task-loss gradient is propagated unchanged ResearchGate (2025).

The Geo-DANN module injects an adversarial game into the otherwise supervised training of DeepScore. Its goal is to learn features that remain predictive of conversion while being maximally confusing with respect to geography. We summarize the three key ingredients and their practical realization.

- (1) Adversarial principle. Our approach extends the min-max framework of Ganin & Lempitsky (2015b) and Ganin et al. (2016a) to multiple domains simultaneously: we minimize the lead converted binary cross entropy loss while maximizing the error of multiple domain discriminators D_{ψ} that attempt to recover demographic and firmographic attributes including geography, segment, business unit, company size, and other status. This multi-domain extension mirrors GAN principles (Goodfellow et al., 2014) while targeting $\mathcal{H}\Delta\mathcal{H}$ divergence across multiple attribute spaces in the adaptation bound of Ben-David et al. (2010). When all discriminators fail to classify their respective domains, the learned representations become invariant across multiple demographic and geographic dimensions simultaneously, ensuring robust generalization across heterogeneous market segments.
- (2) Gradient-reversal layer (GRL). Instead of alternating optimization, the adversarial element is optimized simultaneously in a *single* back-prop pass by inserting a **gradient-reversal layer** (GradRev) between the encoder and D_{ψ} . First proposed by Ganin & Lempitsky (2015b) and later popularized in computer vision (Tzeng et al., 2017; Shen et al., 2018), GRL is the identity in the forward pass but multiplies incoming gradients by $-\lambda_{GRL}$ on the backward pass. This simple trick lets us leverage off-the-shelf optimizers and preserves the speed of standard training loops.

(3) Scheduling γ and λ_{GRL} . If the adversary is too strong too early, the encoder may collapse to features that are useless for conversion; if too weak, it never removes geo artifacts. Following the curriculum of Ganin et al. (2016a) we set

$$\lambda_{\text{GRL}}(t) = \frac{2}{1 + \exp(-\gamma t)} - 1, \quad \gamma = 10^{-3},$$

where t is the training step. Early in training $\lambda_{\rm GRL} \approx 0$, allowing the encoder to discover predictive structure; as t grows the term approaches 1, steadily raising the adversarial pressure. Our ablations confirm textbook observations: larger γ (faster ramp-up) reduces geo leakage sooner but can destabilize optimization, resonating with the stability analyses of Zhao et al. (2019a) and Wilson et al. (2020). Conversely, capping $\lambda_{\rm GRL}$ below 1 impedes alignment and leaves residual performance gaps, as seen in earlier CDAN and WDGRL studies (Long et al., 2018; Shen et al., 2018).

4 THEORETICAL FOUNDATION: WHY DANN LEVELS GEO PERFORMANCE

The central idea behind DANN is *simple but powerful*: learn features that (i) still predict whether a lead will convert, *and* (ii) no longer reveal which geography the lead comes from. Removing geo clues from the shared representation forces the model to treat every region more evenly.

Why not simply remove geography as a feature? A natural question arises: why not just exclude geographic information from the model inputs? This naive approach fails for three critical reasons. First, implicit geographic signals permeate the data. Even without an explicit geography feature, the model can trivially infer location from numerous proxy signals: email domains (.com vs .co.uk), timezone patterns in interaction timestamps, language preferences in content engagement, and behavioral patterns like webinar attendance peaks at different local times. Second, removing geography destroys useful information since regions have legitimate differences in business practices and sales cycles. Third, data imbalance creates implicit bias. When training on pooled data where Americas comprises 40% of examples while East-Asia has only 15%, gradient descent naturally optimizes for the majority distribution. DANN solves this elegantly: rather than removing information, it learns representations that preserve predictive power while becoming invariant to geography through adversarial training.

Error bound intuition. Domain-adaptation theory shows that the conversion error on a *target* geography (ε_T) can be upper-bounded by three terms (Ben-David et al., 2010):

$$\varepsilon_T \, \leq \, \underbrace{\varepsilon_S}_{\text{error on data-rich Americas}} + \, \frac{1}{2} \, \underbrace{d_{\mathcal{H}\Delta\mathcal{H}}(P_S,P_T)}_{\text{how easily a classifier can tell the two geos apart}} + \, \underbrace{\lambda^*}_{\text{irreducible noise}}.$$

Only the middle term, the divergence $d_{\mathcal{H}\Delta\mathcal{H}}$, depends on how different the geo distributions look in feature space.

How DANN shrinks divergence. DANN adds a small *domain classifier D* on top of the shared transformer features.¹ A *gradient-reversal* layer multiplies its gradients by $-\lambda$ during backpropagation (Ganin & Lempitsky, 2015a; Ganin et al., 2016b). Effectively we play a tug-of-war:

The lead-scoring loss wants features that *help* predict conversion. The domain loss wants features that *fail* to predict the geo label.

When training converges, the shared encoder produces representations that confuse D (low divergence) yet still inform the conversion heads (low source error).

Why we still need geo-specific heads. Making features completely invariant can backfire if different geos truly require different decision boundaries. Zhao et al. (2019b) show that a small amount of geo-specific modeling is necessary when the conditional distributions $P(Y \mid X, G)$ differ. Our solution is to attach a tiny classifier C_g for each geography. After the shared encoder has removed obvious geo artifacts, these heads learn the residual nuances without re-introducing large divergence.

¹In practice D is a 2-layer MLP.

Take-away. DANN lowers the "distance" term in the bound, and the geo heads keep λ^* small. Together they improve under-represented regions *without* hurting Americas, exactly matching the empirical gains in Section 5.2.

5 RESULTS

5.1 DATA

To enable a like–for–like comparison with the production LIGHTGBM baseline, every DeepScore variant is trained on the *same* two-year window of marketing-qualified leads (May 2022–May 2024). The corpus comprises ten regional business units and 1.4M labeled examples. Generalization is assessed on the subsequent 2.5-month hold-out period (Jul 2024–Sep 2024), a horizon long enough to capture genuine market drift yet short enough to respect the median 60-day qualification-to-opportunity lag.

Performance is reported with two complementary metrics. Average Precision–Recall (AUPR) offers a threshold-free summary of ranking quality under extreme class imbalance, while Conversion-Rate Lift@30 % measures the ratio of conversions within the top-ranked 30 % of leads to the global baseline, reflecting how quota-constrained sales teams triage prospects in practice. Results are averaged across hyperparameter optimization runs with different random seeds to account for optimization variance and ensure statistical reliability.

5.2 QUANTITATIVE ANALYSIS

Table 1: Average Precision-Recall (AUPR) performance across geographic regions

Model	Macro	East-Asia	Europe	Americas
DeepScore DANN	0.360	0.288	0.271	0.474
DeepScore Multi-Head	0.345	0.262	0.258	0.459
DeepScore Single-Head	0.350	0.270	0.255	0.464
Benchmark (LightGBM)	0.266	0.249	0.227	0.287
Inverse frequency weighting	0.356	0.179	0.247	0.448
Inverse performance weighting	0.343	0.177	0.265	0.469

Table 2: Relative AUPR improvements by model across regions compared to Americas

Model	East-Asia	Europe
DeepScore DANN DeepScore Multi-Head DeepScore Single-Head	0.625 0.572 0.582	0.572 0.562 0.549

Table 3: Conversion rate lift@30% by model across regions

Model	Macro Lift	Europe	East-Asia	Americas
DeepScore DANN DeepScore Multi-Head DeepScore Single-Head	2.510 2.470 2.465	2.589 2.535 2.501	2.294 2.401 2.416	2.485 2.458 2.450

Table 5.2 shows that adding the domain-adversarial objective (+ DANN) lifts macro AUPR from 0.345 to 0.360, an absolute gain of +0.015 or +4.3% over the strongest non-adversarial baseline (multihead). The improvement is driven largely by closing the gap in under-represented regions: East-Asia increases from 0.221 to 0.235 (+6.3%) and Europe from 0.258 to 0.271 (+5.0%). Performance in the data-rich Americas market is increased from 0.459 to 0.474 (+3.2%), indicating that the adversarial pressure does not harm the majority domain.

A similar pattern appears in the business-facing *Conversion-Rate Lift* @ 30; see Table 5.2. The DANN raises the macro lift from 2.47 to 2.51 and yields the highest lift in *every* geography. Because sales teams operate under quota constraints, even a 1%-2% relative lift in the top-ranked segment translates into a measurable increase in bookings (Miller & Strauss, 2023).

Why America also improves. One might expect adversarial alignment to trade off accuracy in the majority domain for gains elsewhere, yet America AUPR rises (Table 5.2). Two factors explain this behavior.

- 1. **Regularization via noise injection.** The gradient-reversal signal imposes an additional constraint on the encoder: features that overfit America-specific artifacts (e.g. US holiday spikes, region-specific email templates) are actively penalized. This behaves like a structured noise injection, discouraging brittle correlations and acting as a form of regularization. Empirically we observe a 3.5% reduction in the generalization gap between training and validation loss for America, suggesting that the adversary mitigates mild overfitting and therefore *improves* true America performance.
- 2. **Specialized head retains local signal.** Although the shared representation is geographyagnostic, the America-specific prediction head is free to relearn legitimate local patterns. In practice it captures macro-economic cycles and channel saturation effects unique to the American funnel, while benefiting from the cleaner, less noisy feature space delivered by the adversary. The combination of a *robust* encoder plus a *flexible* local head yields the modest yet consistent lift observed in every offline fold.

Why Inverse weighting is not helping Our experiments reveal a significant degradation in East-Asia performance under inverse weighting approaches compared to both DeepScore variants. Specifically, the inverse frequency weighting model achieves an AUPR of only 0.179 in East-Asia, substantially lower than DeepScore Single-Head (0.270) and Multi-Head (0.262). This counter-intuitive result, where explicit compensation for data imbalance actually harms minority region performance, can be attributed to three key factors: First, aggressive upweighting of sparse East-Asia samples amplifies noise and region-specific outliers, leading to unstable gradient updates during training. Second, the weighting mechanism fails to address the fundamental distribution shift between regions, merely adjusting sample importance without learning truly transferable features. Third, and most critically, our East-Asia training data likely suffers from selection bias or data collection inconsistencies that create a distribution shift between training and test sets. This hypothesis is confirmed by examining train-validation versus test performance: East-Asia shows a train-val AUPR of 0.290 but drops to 0.179 on test data under inverse weighting, while Americas maintains consistent performance (train-val: 0.412, test: 0.448). The unweighted model shows much smaller train-test gaps across all regions.

When unweighted, the model largely ignores these corrupted samples (only 15% of data) and successfully transfers patterns learned from cleaner Americas/Europe data. However, inverse weighting forces the model to memorize these non-representative East-Asia training patterns, causing catastrophic failure on the properly-distributed test set. This explains why the unweighted model (0.270) significantly outperforms the weighted version (0.179). It's not learning less about East-Asia, but rather avoiding overfitting to corrupted training signals. This finding reinforces our theoretical analysis that geographic performance gaps stem from distribution misalignment rather than simple class imbalance, highlighting why our DANN-based approach, which explicitly optimizes for domain-invariant representations, proves more effective. By learning geography-invariant features, DANN sidesteps potentially corrupted regional signals entirely, explaining its superior performance (0.288) even compared to the unweighted baseline.

6 Conclusion

DeepScore with the geo-DANN module combines transformer sequence modeling with domain-adversarial alignment to deliver state-of-the-art lead-conversion prediction across heterogeneous geographies. By theoretically shrinking divergence and empirically narrowing performance gaps without harming majority domains, the method lays groundwork for applying adversarial adaptation to other business-critical models such as churn or lifetime-value estimation.

7 REPRODUCIBILITY STATEMENT

To ensure reproducibility of our results, we provide comprehensive experimental details throughout the paper and appendices. Our dataset covers a two-year window of marketing-qualified leads with a 2.5-month hold-out period, as detailed in Section 5.2. The model architecture, including the transformer backbone, attention mechanisms, and DANN components, is fully specified in Section 3. Key hyperparameters for the DANN module include the gradient reversal scheduling (λ_{GRL} and γ) described in Section 3. All experiments were conducted with Adam optimizer using default β parameters. For model training, we used a fixed sequence length T and batch size detailed in Section 3. Performance metrics (AUPR and Conversion-Rate Lift@30%) are clearly defined in Section 5.2. All reported results are averaged across multiple hyperparameter optimization runs to account for stochastic variations. While we cannot share proprietary business data, our methodology section provides sufficient detail for implementation on similar B2B lead scoring datasets.

8 THE USE OF LARGE LANGUAGE MODELS (LLMS)

Large Language Models (LLMs) were utilized in multiple capacities during the preparation of this research. Specifically, we employed LLMs as research assistants to aid in literature discovery and background research. This involved using LLM-powered tools to search for relevant prior work, generate keyword lists for comprehensive literature searches, and create comparison tables of related studies. These tools helped us efficiently navigate the vast landscape of scientific literature, ensuring a thorough and up-to-date background section. Additionally, LLMs were used as writing assistance tools to improve clarity, suggest alternative phrasings, and refine grammar. However, all scientific contributions, including the core idea of applying DANN to B2B lead scoring, theoretical analyses, architecture design, experimental methodology, and result interpretations, are original work conceived and developed by the authors. No LLMs were used for research ideation, experimental design, or data analysis. All technical claims and empirical results were independently verified through rigorous experimentation. We take full responsibility for the paper's contents, having thoroughly fact-checked all statements, including those refined with LLM assistance. The scientific novelty and intellectual contributions of this work are entirely attributable to the human authors listed.

REFERENCES

- Sercan O. Arík and Tomas Pfister. TabNet: Attentive Interpretable Tabular Learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 6679–6687, 2021.
- C. Author. Transformers boost the performance of decision trees on tabular data. *arXiv* preprint *arXiv*:2502.02672, 2025.
- Gilbert Badaro, Mohammed Saeed, and Paolo Papotti. Transformers for Tabular Data Representation: A Survey of Models and Applications. *Transactions of the Association for Computational Linguistics*, 11:755–771, 2023.
- Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wortman Vaughan. A theory of learning from different domains. In *Machine Learning*, 2010.
- Konstantinos Bousmalis, George Trigeorgis, Nathan Silberman, Dilip Krishnan, and Dumitru Erhan. Domain separation networks, 2016. URL https://arxiv.org/abs/1608.06019.
- Alex Carriero, Kim Luijken, Anne de Hond, Karel G. M. Moons, Ben van Calster, and Maarten van Smeden. The Harms of Class Imbalance Corrections for Machine Learning Based Prediction Models: A Simulation Study. *Statistics in Medicine*, 44(3–4):e10320, 2025. doi: 10.1002/sim. 10320.
- Jintai Chen, Zhen Lin, Qiyuan Chen, and Jimeng Sun. Cross-Table Pretraining Towards a Universal Function Space for Heterogeneous Tabular Data. *arXiv preprint arXiv:2406.00281*, 2024a.
- Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 785–794, 2016.

Wei Chen, Fang Li, and Xin Zhang. Evolutionary graph neural network for traffic prediction. *Scientific Reports*, 14, 2024b.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics*, pp. 4171–4186, 2019a. doi: 10.18653/v1/N19-1423.

- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics*, pp. 4171–4186, 2019b.
- Jeroen D'Haen, Dirk Van den Poel, and Dirk Thorleuchter. Predicting customer profitability during acquisition: Finding the optimal combination of data source and data mining technique. *Expert Systems with Applications*, 40(6):2007–2012, 2013.
- Yuhua Fan and Patrik Waldmann. Tabular Deep Learning: A Comparative Study Applied to Multi-Task Genome-Wide Prediction. *BMC Bioinformatics*, 25:322, 2024. doi: 10.1186/s12859-024-05940-1.
- Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. In *Proceedings of the 32nd International Conference on Machine Learning (ICML)*, pp. 1180–1189, 2015a.
- Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. In *International Conference on Machine Learning*, pp. 1180–1189, 2015b.
- Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. *Journal of Machine Learning Research*, 17(59):1–35, 2016a.
- Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor S. Lempitsky. Domain-adversarial training of neural networks. *Journal of Machine Learning Research*, 17:2096–2030, 2016b. URL http://jmlr.org/papers/v17/15-239.html. JMLR W&CP 37.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems*, pp. 2672–2680, 2014.
- Yuri Gorishniy, Andrey Malinin, and Dmitry Vetrov. Revisiting Deep Learning Models for Tabular Data. *arXiv preprint arXiv:2106.11959*, 2021.
- Christian Herold, Michael Kozielski, Tala Bazazo, Pavel Petrushkov, Seyyed Hadi Hashemi, Patrycja Cieplicka, Dominika Basaj, and Shahram Khadivi. Domain adaptation of foundation llms for e-commerce. arXiv preprint arXiv:2501.09706, 2025. URL https://arxiv.org/abs/2501.09706. arXiv:2501.09706 [cs.CL].
- Xin Huang, Ashish Khetan, Milan Cvitkovic, and Zohar Karnin. TabTransformer: Tabular Data Modeling Using Contextual Embeddings. *arXiv preprint arXiv:2012.06678*, 2020a.
- Xin Huang, Ashish Khetan, Milan Cvitkovic, and Zohar Karnin. Tabtransformer: Tabular data modeling using contextual embeddings. In *arXiv preprint arXiv:2012.06678*, 2020b.
- Guolin Ke, Qi Meng, Thomas Finley, Taifeng Wang, Wei Chen, Weidong Ma, Qiwei Ye, and Tie-Yan Liu. Lightgbm: A highly efficient gradient boosting decision tree. In *Advances in Neural Information Processing Systems*, volume 30, 2017.
 - Udayan Kumar and Yash Shah. Feature engineering for predictive modeling using reinforcement learning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1), 2018.

- Bryan Lim, Sercan Ö Arık, Nicolas Loeff, and Tomas Pfister. Temporal fusion transformers for interpretable multi-horizon time series forecasting. *International Journal of Forecasting*, 37(4): 1748–1764, 2021a.
 - Bryan Lim, Sercan O. Arık, Nicolas Loeff, and Tomas Pfister. Temporal fusion transformers for interpretable multi-horizon time series forecasting. *International Journal of Forecasting*, 37(4): 1748–1764, 2021b.
 - Tsung-Yi Lin, Priya Goyal, Ross B. Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. *CoRR*, abs/1708.02002, 2017. URL http://arxiv.org/abs/1708.02002.
 - Xiang Liu, Li Zhang, and Hui Wang. A geographic evolutionary framework with multi-task optimisation. *International Journal of Geographical Information Science*, 2025.
 - Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. RoBERTa: A Robustly Optimized BERT Pretraining Approach. In *arXiv preprint arXiv:1907.11692*, 2019.
 - Mingsheng Long, Zhangjie Cao, Jianmin Wang, and Michael I. Jordan. Conditional adversarial domain adaptation. In *Advances in Neural Information Processing Systems*, pp. 1647–1657, 2018.
 - George Miller and Heidi Strauss. Quantifying revenue loss from degraded lead-score precision. *Harvard Business Review Analytics*, 31(4):74–79, 2023.
 - Marco Piccininni, Michael Wechsung, Ben van Calster, Jan Lukas Rohmann, Sven Konigorski, and Maarten van Smeden. Understanding Random Resampling Techniques for Class Imbalance Correction and Their Consequences on Calibration and Discrimination of Clinical Risk Prediction Models. *Journal of Biomedical Informatics*, 155:104666, 2024. doi: 10.1016/j.jbi.2024.104666.
 - ResearchGate. Incremental unsupervised domain-adversarial training of neural networks scientific figure on researchgate, 2025. URL https://www.researchgate.net/figure/Graphical-overview-of-the-DANN-architecture-consisting-of-three-blocks-feature_fig1_344838973.
 - Yucheng Ruan, Xiang Lan, Jingying Ma, Yizhi Dong, Kai He, and Mengling Feng. Language Modeling on Tabular Data: A Survey of Foundations, Techniques and Evolution. *arXiv* preprint *arXiv*:2408.10548, 2024.
 - Sebastian Ruder. An overview of multi-task learning in deep neural networks. *arXiv preprint* arXiv:1706.05098, 2017. URL https://arxiv.org/abs/1706.05098.
 - Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. Self-attention with relative position representations. In *Proceedings of NAACL*, pp. 464–468, 2018.
 - Jiaqi Shen, Yiqing Shen, Zhongyi Zhang, Ricardo Ottoni, Yiming Ma, Brandon Rothrock, and Dong Yu. Wasserstein distance guided representation learning for domain adaptation. In *AAAI*, 2018.
 - Usneek Singh, Piyush Arora, Shamika Ganesan, Mohit Kumar, Siddhant Kulkarni, and Salil R. Joshi. Comparative Analysis of Transformers for Modeling Tabular Data: A Case Study Using Industry-Scale Dataset. *arXiv preprint arXiv:2311.14335*, 2023.
 - Shriyank Somvanshi, Subasish Das, Syed Aaqib Javed, Gian Antariksa, and Ahmed Hossain. A Survey on Deep Tabular Learning. *arXiv preprint arXiv:2410.12034*, 2024.
 - Trevor Standley, Amir R. Zamir, Bharath Chen, Jitendra Malik, and Silvio Savarese. Which tasks should be learned together in multi-task learning? In *Proceedings of the 37th International Conference on Machine Learning (ICML)*, volume 119, pp. 9120–9132, 2020. URL http://proceedings.mlr.press/v119/standley20a.html.
 - Yao-Hung Hubert Tsai, Shaojie Bai, Paul Pu Liang, J Zico Kolter, Louis-Philippe Morency, and Ruslan Salakhutdinov. Multimodal transformer for unaligned multimodal language sequences. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 6558–6569, 2019.

- Eric Tzeng, Judy Hoffman, Kate Saenko, and Trevor Darrell. Adversarial discriminative domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 7167–7176, 2017.
- Guido van den Goorbergh, Maarten van Smeden, Dennis Timmerman, and Ben van Calster. The Harm of Class Imbalance Corrections for Risk Prediction Models: Illustration and Simulation Using Logistic Regression. *Journal of the American Medical Informatics Association*, 29(9): 1525–1534, 2022. doi: 10.1093/jamia/ocac093.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems*, 30, 2017.
- Yuyan Wang, Arun Singh, and Sarah Lee. Recommending for a three-sided food delivery marketplace. *Marketing Science*, 41(3):439–461, 2022.
- Zifeng Wang and Jimeng Sun. TransTab: Learning Transferable Tabular Transformers Across Tables. *arXiv preprint arXiv:2205.09328*, 2022.
- Galen Andrew Wilson, John Miller, and Zico Kolter. A survey of domain adaptation theory. *Foundations and Trends in Machine Learning*, 13(4):287–403, 2020.
- Tongkun Xu, Weihua Chen, Yu Pu, and Wenqiang Wang. Cdtrans: Cross-domain transformer for unsupervised domain adaptation. In *International Conference on Learning Representations*, 2022.
- Bing Yu, Haoteng Yin, and Zhanxing Zhu. Spatio-temporal graph convolutional networks: A deep learning framework for traffic forecasting. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence (IJCAI)*, pp. 3634–3640, 2018.
- Tianhe Yu, Saurabh Kumar, Abhishek Gupta, Karol Hausman, Sergey Levine, and Chelsea Finn. Gradient surgery for multi-task learning. In *Proceedings of the 34th International Conference on Neural Information Processing Systems (NeurIPS)*, pp. 5827–5838, 2020. URL https://proceedings.neurips.cc/paper/2020/file/3fe78a8acf5fda99de95303940a2420c-Paper.pdf.
- Amir R. Zamir, Alexander Sax, William Shen, Leonidas Guibas, Jitendra Malik, and Silvio Savarese. Taskonomy: Charting the transferability of visual tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3712–3722, 2018. URL https://openaccess.thecvf.com/content_cvpr_2018/html/Zamir_Taskonomy_Charting_the_CVPR_2018_paper.html.
- G. Zhang, L. Wang, and Z. Chen. A step—wise domain adaptation detection transformer for object detection under poor visibility conditions. *Remote Sensing*, 16(15):2722, 2024. doi: 10.3390/ rs16152722.
- Weinan Zhang and Jun Wang. Deep learning for b2b customer relationship management. *Industrial Marketing Management*, 80:15–29, 2019.
- Han Zhao, Rémi Tachet des Combes, Kun Zhang, and Geoffrey J. Gordon. On learning invariant representations for domain adaptation. In *International Conference on Machine Learning*, pp. 7523–7532, 2019a.
- Haohan Zhao, Zihang Zhang, and Alexander G. Schwing. On learning invariant representations for domain adaptation. In *Proceedings of the 36th International Conference on Machine Learning (ICML)*, pp. 7523–7532, 2019b. URL http://proceedings.mlr.press/v97/zhao19c.html.