# TATN-Tri-Modal Adversarial Transparency Network (TATN)

---

## 1   Introduction

TATN(Tri-modal Adversarial Transparency Network) combines:
1. **DSCD** — Dynamic Span–Sense Co–Detection: an online prototype-based module that jointly detects ambiguous spans (homographs) and predicts sense cluster IDs while computing multi-source uncertainty; and
2. **ASBN** — Adversarial Sense Balance Network: training-time discriminators that detect shortcut behaviors (frequency bias, context-sparsity guessing, and cross-lingual mismatch) and push the encoder to remove them using a Gradient Reversal Layer (GRL).

The following sections present the algorithmic details in an Overleaf-friendly format: short explanations followed by compact pseudocode suitable for inclusion in a paper.

## 2   Notation (quick reference)

$S$ — Source sentence (subwords) $S = [s_1, \ldots, s_L]$.

$w_j$ — Token type (subword string) at position $j$.

$h_j \in \mathbb{R}^d$ — Contextual encoder vector for token $j$.

$C_w = \{c_{w,1}, \ldots, c_{w,K_w}\}$ — Dynamic prototype set for type $w$.

$p_j$ — Soft distribution over prototypes for token $j$.

$\hat{y}_j$ — Selected prototype / sense id $(p_j)$.

$p_{\max}(j)$ — Confidence = $\max_k p_{j,k}$.

$U_j$ — Combined uncertainty score (entropy + MC var + learned $\sigma$ + novelty).

$g_j$ — Uncertainty gate (attention boost scalar in $(0, 1)$).

$z_{a(j)}^{\mathbf{tgt}}$ — Aligned target embedding (parallel corpora).

$D_{\mathbf{freq}}, D_{\mathbf{ctx}}, D_{\mathbf{xl}}$ — ASBN discriminators.

$\bar{\lambda}_k$ — base GRL strength for discriminator $k$.

# 3   High-level pipeline (what/why)

1. **DSCD forward pass (inference  training):** encode $\rightarrow$ compute multi-source uncertainty $\rightarrow$ dynamic prototype assignment/create/update $\rightarrow$ attention gating $\rightarrow$ joint span detection + sense assignment $\rightarrow$ sense-augmented token vectors $\{h'_j\}$ $\rightarrow$ decoder.

2. **ASBN (training only):** three lightweight discriminators monitor DSCD outputs and send reversed gradients (via GRL) to the encoder to eliminate shortcut signals. ASBN uses per-sample confidence weighting to avoid penalizing genuine uncertainty.

3. **TRG (optional):** rationale generator that consumes DSCD+ASBN diagnostics to create human-readable explanations for decisions.

# 4   Core DSCD forward-pass (compact pseudocode)

Below is a compact algorithm for the DSCD forward pass. The comments explain "what" is computed and "why."

---

**Algorithm 1** DSCD forward pass (per sentence)

---

**Require:** Source sentence $S = [s_1, \ldots, s_L]$

1: **for** $j = 1$ to $L$ **do**                      $\triangleright$ Encode and buffer
2:      $h_j \leftarrow \text{Encoder}(s_j)$                 $\triangleright$ contextual embedding
3:      append $h_j$ to buffer $B_{w_j}$          $\triangleright$ maintain recent embeddings per type
4: **end for**
5: Compute per-type dispersion $D_w$ using buffers $B_w$      $\triangleright$ detect multi-sense candidates
6: **for** $j = 1$ to $L$ **do**             $\triangleright$ Per-token uncertainty and prototype logic
7:      **if** $C_{w_j} = \emptyset$ **then**
8:          create initial prototype (immediately or after $N_{\min}$ buffered examples)
9:      **else**
10:          compute cosine similarities $s_{j,i} = \cos(h_j, c_{w_j,i})$ for $i = 1..K_{w_j}$
11:          $p_j \leftarrow \text{softmax}(s_{j,\cdot}/T)$          $\triangleright$ soft assignment to prototypes
12:          $H_j \leftarrow -\sum_i p_{j,i} \log(p_{j,i} + \varepsilon)$          $\triangleright$ aleatoric entropy
13:          compute MC-dropout variance $\text{Var}_j$ (run $M$ stochastic passes)
14:          compute learned noise $\sigma_j$ from a small $\sigma$-Net
15:          $d_{\min} \leftarrow \min_i(1 - s_{j,i})$          $\triangleright$ novelty to nearest centroid
16:          $U_j \leftarrow \alpha_1 H_j + \alpha_2 \sigma_j + \alpha_3 \text{Var}_j + \alpha_4 \text{norm}(d_{\min})$
17:      **end if**
18:      Compute adaptive threshold $\varepsilon_{w_j}^{new}$ from rolling mean/std of assignment distances
19:      **if** $d_{\min} > \varepsilon_{w_j}^{new}$ **and** $U_j > \delta_{\text{inst}}$ **then**
20:          create new prototype $c_{w_j,K+1} \leftarrow h_j$
21:      **else**
22:          assign token to nearest prototype $i^\star \leftarrow_i s_{j,i}$
23:          update centroid $c_{w_j,i^\star} \leftarrow (1 - \eta)c_{w_j,i^\star} + \eta\, h_j$
24:      **end if**
25: **end for**
26: For flagged tokens (types with $D_w > \delta_{\text{type}}$ and $U_j > \delta_{\text{inst}}$) compute gate $g_j = \sigma(w_g(U_j - b_g))$
27: Boost base attention $a_j^{(0)} \leftarrow a_j^{(0)}(1 + \gamma g_j)$ and renormalize to $\tilde{a}_j$
28: Span prediction $\hat{b}_j \leftarrow \sigma(W_{\text{span}} h_j + b_{\text{span}})$ (binary)
29: Sense id $\hat{y}_j \leftarrow_i p_{j,i}$ (or clustering assignment)
30: Form sense-augmented vector:

$$h'_j \leftarrow \begin{cases} h_j + c_{w_j,\hat{y}_j} & \text{if } \hat{b}_j > 0.5 \\ h_j & \text{otherwise} \end{cases}$$

31: Send $\{h'_j\}, \{\tilde{a}_j\}$ to decoder

---

**Explanation (short):**

- The forward pass computes robust, multi-source uncertainty signals so the system knows which tokens are ambiguous.

- Prototypes are created/updated online using adaptive (per-type) thresholds, and EMA updates keep centroids stable.

- Flagged tokens receive larger attention via a learnable gate so the decoder sees more context about ambiguous positions.

# 5   ASBN — training-time adversarial module (plain description)

ASBN consists of three lightweight discriminators:

$D_{\mathbf{freq}}$ Frequency shortcut detector: checks whether the model's sense choice equals the most frequent sense for the token type.

$D_{\mathbf{ctx}}$ Context-sparsity detector: checks whether the encoder is ignoring weak contexts and defaulting to priors.

$D_{\mathbf{xl}}$ Cross-lingual detector: checks whether the source predicted sense aligns with the target token embedding (uses negatives created from wrong-sense prototypes).

Each discriminator receives the encoder representation *through a GRL*. The GRL forwards identity in the forward pass and multiplies gradients by $-\lambda$ in the backward pass, causing the encoder to learn representations that *hide* the shortcuts the discriminators exploit.

Per-sample GRL strengths are computed to avoid punishing honest uncertainty:

$$\lambda_{k,j} = \mathrm{clip}\big(\bar{\lambda}_k \cdot p_{\max}(j) \cdot (1 - U_j) \cdot g_j, \ 0, \ \lambda_{\max}\big)$$

where $k \in \{\mathrm{freq}, \mathrm{ctx}, \mathrm{xl}\}$.

# 6   ASBN pseudocode and integrated training step

The pseudocode below shows a single training-step (per batch) that integrates DSCD forward-pass (Algorithm 1) with ASBN discriminator updates and GRL-based adversarial training.

---

**Algorithm 2** TATN v2 training step (per batch)

---

**Require:** batch of parallel sentences (source, target)

1: Run DSCD forward for each source sentence $\Rightarrow$ obtain $\{h_j, p_j, \hat{y}_j, p_{\max}(j), U_j, g_j\}$ and update prototypes (Alg. 1)

2: Compute decoder loss $\mathcal{L}_{\text{MT}}$ (translation CE), span loss $\mathcal{L}_{\text{span}}$ and sense loss $\mathcal{L}_{\text{sense}}$

3: **for** each token $j$ in batch **do**

4:     Compute per-discriminator GRL weights:

$$\lambda_{k,j} \leftarrow \text{clip}\big(\bar{\lambda}_k \cdot p_{\max}(j) \cdot (1 - U_j) \cdot g_j, \ 0, \lambda_{\max}\big)$$

5: **end for**

6:                              $\triangleright$ — Discriminator forward passes (via per-sample GRL) —

7: $x_j^{\text{freq}} \leftarrow \text{GRL}(h_j; \lambda_{\text{freq},j}) \oplus F_{w_j}$ and compute $\ell_j^{\text{freq}} \leftarrow D_{\text{freq}}(x_j^{\text{freq}})$

8: $x_j^{\text{ctx}} \leftarrow \text{GRL}(h_j; \lambda_{\text{ctx},j}) \oplus \text{ctx\_stats}_j$ and compute $\ell_j^{\text{ctx}} \leftarrow D_{\text{ctx}}(x_j^{\text{ctx}})$

9: $x_j^{\text{xl}} \leftarrow \text{GRL}(h_j; \lambda_{\text{xl},j}) \oplus \phi_{\text{tgt}}(z_{a(j)}^{\text{tgt}})$ and compute contrastive loss $\ell_j^{\text{xl}} \leftarrow D_{\text{xl}}(x_j^{\text{xl}})$

10: Aggregate discriminator losses:

$$\mathcal{L}_{\text{ASBN}} \leftarrow \frac{1}{N} \sum_j \big(w_{\text{freq}}\ell_j^{\text{freq}} + w_{\text{ctx}}\ell_j^{\text{ctx}} + w_{\text{xl}}\ell_j^{\text{xl}}\big)$$

11:                              $\triangleright$ Update discriminator parameters $\phi$ (minimize $\mathcal{L}_{\text{ASBN}}$)

12: `opt_phi.zero_grad();  L_ASBN.backward(retain_graph=True);  opt_phi.step();`

13:                              $\triangleright$ Update main model parameters $\theta$ (encoder, DSCD, decoder)

14: `opt_theta.zero_grad();`

15: $\mathcal{L}_{\text{total}} \leftarrow \mathcal{L}_{\text{MT}} + \lambda_{\text{span}}\mathcal{L}_{\text{span}} + \lambda_{\text{sense}}\mathcal{L}_{\text{sense}} + \lambda_{\text{ASBN}}\mathcal{L}_{\text{ASBN}} + \lambda_{\text{reg}}\mathcal{R}$

16: `L_total.backward();  opt_theta.step();`

17: **return** updated parameters, logs (discriminator accuracies, homograph metrics, prototype stats)

---

**Key remarks (practical):**

- **GRL behavior:** forward pass identity; backward pass multiplies gradients by $-\lambda_{k,j}$ (frameworks like PyTorch can implement sample-weighted GRL or emulate it by weighting losses).

- **Confidence weighting:** we scale GRL so only *wrong-but-confident* shortcuts receive strong adversarial pressure.

- **Discriminator capacity:** keep discriminators small (1–2 layer MLPs) to avoid overpowering the encoder.

- **Warm-up:** train DSCD+NMT for 1–3 epochs before enabling full ASBN to let uncertainty estimates become meaningful.

# 7  Hyperparameters and monitoring

**Suggested ranges**

- Buffer size $C$: 100–1000; EMA $\eta$: 0.03–0.08.

- MC-dropout passes $M$: 3–8; temperature $T$: 0.6–1.0.

- $N_{\min}$ (prototype stabilization): 3–10; $K_{\max}$: 20.

- GRL base strengths: $\bar{\lambda}_{\text{freq}} = 1.0$, $\bar{\lambda}_{\text{ctx}} = 0.5$, $\bar{\lambda}_{\text{xl}} = 0.8$; $\lambda_{\max} = 2.0$.

- Learning rates: $lr_\theta = 3\text{e-}4$, $lr_\phi = 1\text{e-}4$.

**What to log**

- Discriminator accuracy and loss over time (should decrease if encoder hides shortcuts).

- Prototype creation rate, prototype counts per type, centroid drift.

- Homograph-slice translation metrics (Sense-F1, BLEU/COMET on homograph sentences).

- Distribution of $U_j$ (uncertainty) and per-sample $\lambda_{k,j}$.

# 8 Practical tips and common debugging steps

- If discriminators dominate training: reduce $\bar{\lambda}_k$, reduce discriminator size, or decrease discriminator LR.

- If translation quality drops: lower $\lambda_{\text{ASBN}}$ or only apply ASBN to a curated homograph-focused batch.

- If uncertainty $U_j$ is noisy: increase warm-up period and verify MC-dropout/$\sigma$-net implementations.

- Use small, progressive changes: enable one discriminator at a time (first frequency, then context, then cross-lingual).