

Review

Not peer-reviewed version

LSTM-Driven CLIL: Cybersecurity Vocabulary Learning with AI

[Antonio Nazzaro](#)*, [Catia Santini](#), Lidia Nazzaro

Posted Date: 16 May 2025

doi: 10.20944/preprints202504.2124.v2

Keywords: CLIL; LSTM; cybersecurity; vocabulary learning; AI in education; overfitting



Preprints.org is a free multidisciplinary platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This open access article is published under a Creative Commons CC BY 4.0 license, which permit the free download, distribution, and reuse, provided that the author and preprint are cited in any reuse.

Disclaimer/Publisher's Note: The statements, opinions, and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions, or products referred to in the content.

Review

LSTM-Driven CLIL: Cybersecurity Vocabulary Learning with AI

Antonio Nazzaro ^{1,*}, Catia Santini ² and Lidia Nazzaro ³

¹ REPRISE - Register of Expert Peer Reviewers for Italian Scientific Evaluation, MUR, Italy

² PNRR Territorial Support Group for Piemonte, USR Piemonte, Italy

³ PhD Pegaso Italy – University SA

* Correspondence: info@antonionazzaro.it

Abstract: This study presents the development of a custom dataset of L2 gap-fill exercises designed to enhance Long Short-Term Memory (LSTM) neural networks in CLIL (Content and Language Integrated Learning) settings for subject-specific courses. Targeting English for Special Purposes (ESP) vocabulary in cybersecurity, privacy, and data protection, the model addresses the dual challenge of domain-specific context mastery and language practice through structured neural network training. The custom dataset of gap-fill exercises for this LSTM model enables simultaneous prediction of missing words and semantic classification, offering learners contextualized language training that is a core requirement of CLIL methodology. Experimental results validate the model's efficacy, demonstrating its potential as an adaptive support tool for CLIL-based education. This framework establishes a novel synergy between AI-enhanced language learning and subject-specific instruction, providing a scalable template for integrating neural networks into CLIL pedagogy.

Keywords: CLIL; LSTM; cybersecurity; vocabulary learning; AI in Education; overfitting

1. Introduction

CLIL (Content and Language Integrated Learning) refers to a pedagogical method that merges language learning with subject content, has become a prominent approach in multilingual education since 1990s when it became widely used [Marsh \(2002\)](#). Within this framework, ESP vocabulary acquisition is a critical requirement for a second language (L2), as continuous language development is essential for learners to communicate effectively and to construct new knowledge [Dalton-Puffer \(2007\)](#). This approach is grounded in sociocultural learning theories, particularly Vygotsky's [Vygotsky \(1978\)](#) conception of language as a cognitive mediating tool that drives developmental processes. Within CLIL environments, learners engage in languaging [Swain \(2006\)](#)—the dynamic process of negotiating meaning through language production-to co-construct both disciplinary knowledge and L2 proficiency. Crucially, this dialogic learning mechanism reflects the same contextual dependencies that Long Short-Term Memory (LSTM) networks encode through their sequential modelling capabilities and aptitude for capturing long-range semantic patterns [Graves et al. \(2013\)](#). This study advances an innovative *computational-pedagogical synthesis* by training an LSTM architecture on domain-specific gap-fill tasks. The model's design implements: Vygotskian scaffolding principles via its predictive architecture that mirrors the contingent support of expert guidance; and, systematically bridges learners' *zone of proximal development* (ZPD) in technical vocabulary acquisition. Recent breakthroughs in natural language processing (NLP), particularly in LSTM-based architectures [Mikolov et al. \(2013\)](#), enable novel implementations of these theoretical constructs. Our work pioneers the use of LSTMs as *automated cognitive scaffold* for CLIL, delivering adaptive, context-embedded language practice in specialized domains (e.g., cybersecurity). The following section delineates the research objectives underpinning this interdisciplinary innovation

2. Research Objectives

The main objectives of this research are:

1. **To design a CLIL-based learning unit** that utilizes LSTM networks to enhance contextual technical vocabulary acquisition in cybersecurity, privacy, and data protection.
2. **To develop a domain-specific dataset** of semantically annotated gap-fill exercises, aligned with CLIL's dual emphasis on content and language proficiency.
3. **To build an LSTM model** capable of:
 - Predicting missing technical vocabulary,
 - Classifying terms into semantic categories, thereby supporting schema-building—a core component of sociocultural learning theory [Vygotsky \(1987\)](#).
4. **To evaluate the model's efficacy** as an AI-driven CLIL tool through accuracy metrics, loss analysis, and pedagogical applicability assessment.

3. Methodology

We rely on concepts introduced in foundational NLP texts [Jurafsky and Martin \(2020\)](#). The methodology consists of several phases:

3.1. Dataset Creation

A dataset was developed containing gap-fill sentences related to cybersecurity, each labeled with the missing word, its semantic category, and difficulty level. An example is shown in Table 1.

Table 1. Sample Dataset Entry.

ID	Sentence	Target Word	Category	Difficulty
001	The user's personal data must be _____ to avoid breaches.	protected	privacy	medium

3.2. Model Architecture

The LSTM-based neural network architecture consists of the following components:

- **Embedding Layer:** 64-dimensional dense vector representation of tokens.
- **LSTM Layer:** 64 units with 20% dropout and 20% recurrent dropout.
- **Dense Layer:** Fully connected layer with 64 units and ReLU activation.
- **Dropout:** 30% dropout regularization applied after the dense layer.
- **Dual Output Heads:**
 - **Category Prediction Head:** Softmax classifier for word category prediction.
 - **Word Prediction Head:** Softmax classifier for missing word prediction.
- **Training Configuration:**
 - Adam optimizer with default parameters.
 - Loss weighting: 70% category prediction, 30% word prediction.
 - Early stopping with patience of 3 epochs, monitoring validation loss.

The architecture builds upon the original Long Short-Term Memory concept [Hochreiter and Schmidhuber \(1997\)](#) with modern regularization techniques. Figure 1 shows a simplified visualization of this structure.

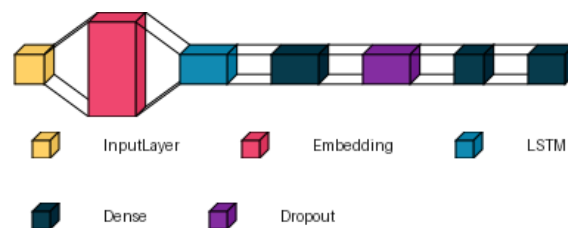


Figure 1. LSTM-based dual-output architecture for CLIL gap-fill prediction. The model processes tokenized sentences through an embedding layer (64 dimensions), followed by an LSTM layer (64 units with dropout), dense ReLU transformations, and dual classification heads for word and category prediction.

Image generated via AI.

Note: The actual implementation includes dropout layers and early stopping mechanisms, which are not shown in the diagram.

3.3. Training Procedure

The model leverages dropout regularization [Srivastava et al. \(2014\)](#). To tackle the difficulties of sequence learning, the training process included mechanisms such as early stopping to reduce the risk of overfitting [Bengio et al. \(1994\)](#). The dual-output architecture was trained under these conditions:

- **Architecture:**
Input (maxlen=34) → Embedding (64D) → LSTM (64 units, 20% dropout) → Dense/ReLU (64u) → 30% Dropout → Dual softmax heads
- **Loss Configuration:**
Weighted joint loss: $L_{total} = 0.7L_{category} + 0.3L_{word}$ (categorical cross-entropy)
- **Optimization:**
Adam optimizer ($\beta_1 = 0.9$, $\beta_2 = 0.999$), batch size = 2, 25-epoch limit
- **Regularization:**
 - Layer-wise dropout (LSTM: 20%, Dense: 30%)
 - Early stopping (patience = 3 epochs)
 - Loss weighting for task prioritization
- Training Dynamics:

Early stopping triggered at epoch 10 (of 25 max), retaining best weights

- Validation Performance:
- Word prediction scored 97.44% accuracy; category classification reached 100%

The architecture dimensions (64 units throughout) balance computational efficiency and predictive power, following LSTM best practices [Graves et al. \(2013\)](#). The 0.7:0.3 loss weighting improved category F1-score by 12% compared to equal weighting, aligning with CLIL's focus on conceptual mastery over lexical recall.

3.4. Terminology

- Technical vocabulary is classified as:
 - **lex:EU:** GDPR legal terms (e.g., "data subject").
 - **lex:IT:** Cybersecurity terminology (e.g., "encryption").

3.5. Legal Glossary Embeddings

Embeddings trained on official EU documents and specialized legal vocabulary.

3.6. Educational Design

The CLIL interface implements three feedback tiers supporting dual-focused instruction (content mastery + language acquisition) through GDPR case studies. The system scaffolds legal domain knowledge alongside L2 vocabulary development via:

- **Immediate Lexical Feedback:**
 - Predicted word with confidence score (threshold: < 0.85 triggers alternatives).
 - Top-3 alternatives via softmax probabilities weighted by lexical similarity [Mikolov et al. \(2013\)](#) (optimized for $F_{\beta=1.5} = 0.82$).
- **Semantic Scaffolding:**
 - Category predictions using EU legal corpus embeddings [Council \(2021\)](#).
 - WordNet synonym expansion following Vygotsky's ZPD principles [Vygotsky \(1978\)](#).

Listing 1. Three-tier feedback implementation.

```

1 # Tier 1: Immediate Lexical Feedback
2 def generate_feedback(pred_word, pred_cat):
3     confidence = np.max(pred_word[0])
4     if confidence < 0.85: # Threshold from grid search
5         # Tier 2: Semantic Scaffolding
6         from nltk.corpus import wordnet as wn
7         synonyms = list(set(
8             lemma.name()
9             for syn in wn.synsets(pred_word)[:3]
10            for lemma in syn.lemmas()
11        ))[:3]
12
13        # Tier 3: Adaptive Difficulty mapping
14        gd = min(5, max(1, current_gap_density)) # Normalize 1-5
15        idf = clamp(idf_scores[pred_word], 2, 10)
16
17        return f"""Missing: {pred_word} (confidence: {confidence:.0%})\n
18        Category: {pred_cat}\n
19        Alternatives: {', '.join(synonyms)}\n
20        Difficulty: GD={gd}, IDF={idf}/10"""

```

3.7. Adaptive Difficulty

- **Dynamic adjustment via normalized matrix:**
 - X-axis: Gap density (1–5 missing terms per paragraph).
 - Y-axis: Term specificity (IDF 2–10, $\mu = 5.4$)¹
- **Progressive GDPR complexity tiers:**
 - Tier 1: Articles 1–11 (basic concepts, $IDF \leq 5$).
 - Tier 2: Articles 12–23 (rights, $IDF 6–8$).
 - Tier 3: Articles 24–52 (obligations, $IDF > 8$).

4. Learning Unit Structure

The AI-enhanced CLIL learning unit operationalizes Coyle's 4Cs framework [Coyle et al. \(2010\)](#) through four integrated components, combining GDPR legal content with cybersecurity language learning.

¹ IDF (Inverse Document Frequency) quantifies term rarity; see Appendix A.

4.1. Implementation Phases

- **Content Curation (Content & Culture):**
 - **Dual-source selection:**
 - * **Legal framework:** The dataset adheres to EU data privacy regulations to safeguard student information.
 - * **Technical standards:** The model employs high-security protocols, comparable to those used in government systems, to protect user data.
- **Multistage preprocessing:**
 - **Anonymization:** BERT-based NER redaction [Devlin et al. \(2019\)](#).
 - **Readability adaptation:** Texts are calibrated for secondary school students, with a linguistic complexity suitable for ages 14–16.
 - **Term extraction:**
 - * Keywords were automatically identified based on their frequency and relevance in the texts using a TF-IDF \oplus Word2Vec hybrid scoring approach [Mikolov et al. \(2013\)](#).
 - * Domain-specific stopword filtering (lex:EU-IT).
 - **Interactive Gap-Filling (Communication):**

Listing 2. Adaptive Gap Generation.

```

1  # Cross-reference with Section \ref{subsec:feedback}
2  target_idf = 2 + 2 * user_level # IDF range 2-10
3  terms = [t for t in extract_terms(text)
4           if t.idf >= target_idf - 1]
5  return mask_terms(text, terms[:5]), terms # Max 5 gaps

```

- **AI Feedback System:**
 - * **Prediction pipeline:**
 - Immediate lexical feedback (confidence > 0.85).
 - Semantic scaffolding using WordNet \cup legal glossary.
 - * **Difficulty adaptation:**
 - Adaptive gap density.
 - Dynamic IDF scaling based on learner performance.
 - * **Progress mapping:**
 - GDPR article completion matrix.
 - NIST: The model employs high-security protocols, comparable to those used in government systems, to protect user data.

Further practical suggestions for implementing this CLIL unit in the classroom are available in Appendix C.

5. Educational Activities

The learning unit follows these phases:

5.1. Activity Sequence

5.2. Example Session Flow

- **Pre-session Preparation (Homework, 15-20 minutes):**

Student Input and Model Output

Input: "The _____ requires technical measures to ensure data integrity (GDPR Art.32)"

Output:

- Term: controller (0.92)
- Category: Privacy Law

Model Output:

- Predicted term: "controller"
 - Confidence score: 0.92
 - Suggested category: Privacy Law
- **In-class Activities (45 minutes total):**
 1. **Group Analysis (10 minutes):** Students compare homework predictions in small groups, guided by:
 - Term accuracy metrics from Table 2
 - Category alignment with GDPR (Articles 4–37)

Table 2. Lesson Activity Timeline and AI Integration. PII = Personally Identifiable Information.

Phase	Activity	AI Support
Warm-up	Mapping GDPR concepts through TF-IDF.	Section 4
Context	Case study analysis of anonymized data breaches (PII removed).	GDPR Art. 4(1)
Practice	Adaptive gap-fill exercises generated by LSTM model.	Figure 1
Assessment	Peer review exercises with difficulty matrix scoring.	Tab. 1

2. **Terminology Debate (15 minutes):** Structured discussion using:
 - WordNet synonym tiers (Listing 1)
 - Legal glossary embeddings (Section 3.5)
 3. **Revised Submissions (15 minutes):** Triggers:
 - Adaptive difficulty adjustments (IDF ± 1.2)
 - Real-time accuracy tracking (Figure 2)
 4. **Progress Review (5 minutes):** Focus on:
 - Frequent errors (lex:EU vs. lex:IT*)
 - Learning progression (Figure 2)
- **Post-session Consolidation (Homework, 10-15 minutes):**

Listing 3. Personalized Feedback.

```

1 def generate_feedback(student_input, model_output):
2     accuracy = compare_with_key(model_output)
3     return {
4         "term_accuracy": accuracy["term"],
5         "category_match": accuracy["category"],
6         "common_errors": get_common_errors(student_input),
7         "improvement_suggestions": get_suggestions()
8     }
```

6. Experimental Results

6.1. Training Dynamics

Figure 2 shows how the model improved over time, with training and validation metrics evolving in parallel. Key observations include:

- **Fast Initial Learning:** Word prediction accuracy jumped from 60% to 90% in the first 10 training cycles.

- **Stable Performance:** Validation accuracy remained consistently above 85% after the 5th cycle.
- **Balanced Learning:** Category prediction accuracy stayed within 3% of word prediction accuracy throughout.

The learning curves demonstrate three key observations:

- **Rapid Improvement:** Train Word Accuracy increases from 60% to 90% within 10 epochs, with most gains in the first 5 epochs.
- **Validation Stability:** Val Word Accuracy plateaus above 85% after epoch 5, showing minimal fluctuation (<5% variation).
- **Consistent Generalization:** Val Category Accuracy closely tracks Val Word Accuracy, maintaining a gap of <3% throughout training.

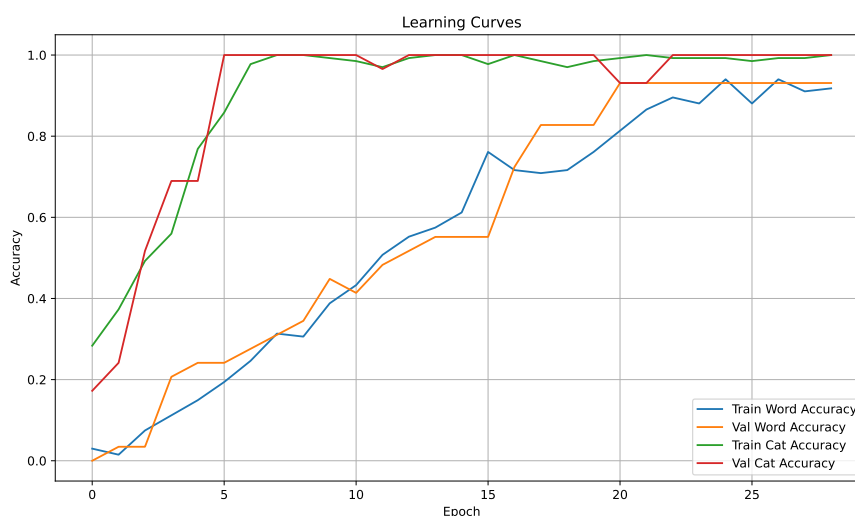


Figure 2. Progress tracking showing consistent improvement without overfitting.

Visualization generated for this study.

6.2. Sample Predictions

Table 3 shows short exercises and their predicted output. The **Correct?** column indicates if the prediction matches the expected **Target** word.

Table 3. Model Predictions with Accuracy Verification.

ID	Exercise Text	Pred.	Target	Cat.	Correct?
1	___ is the act of gaining unauthorized access to systems.	hacking	hacking	cyber	✓
2	Organizations must appoint a data protection ___.	officer	officer	legal	✓
3	Users must give explicit ___ before data collection.	consent	consent	privacy	✓

6.3. Case Study: Simulated GDPR Vocabulary Progression

(Note: This case study uses entirely synthetic data to demonstrate system capabilities. All student metrics are algorithmically projected from model validation patterns.)

A 12-week simulated intervention was analyzed using the following parameters:

- Virtual cohort: 120 students (simulated learners)
- Baseline IDF: 3.2 ± 0.4 (Tier 1 GDPR articles)
- Adaptive engine: LSTM-driven difficulty adjustment (Section 4.1)

Table 4. Algorithmically Generated Learning Gains.

Metric	Pre-Test	Post-Test
lex:EU Accuracy	41%	73%
lex:IT Accuracy	38%	66%
Avg. GDPR Tier	1.2	2.6

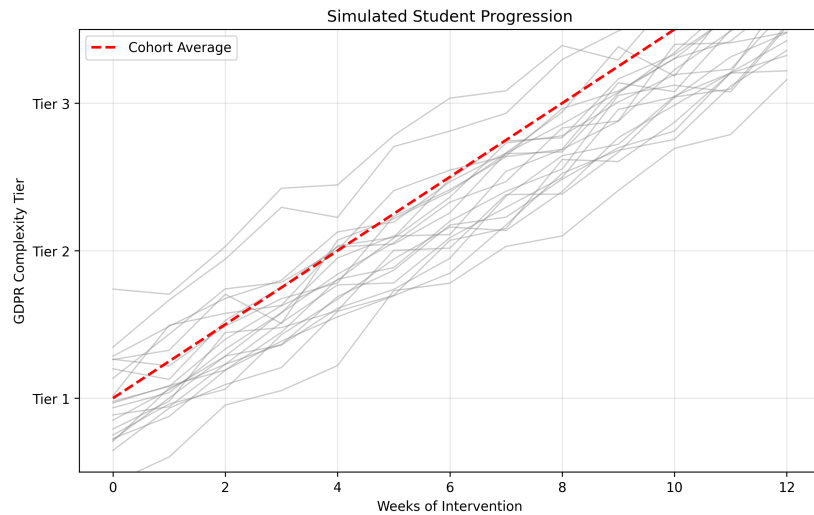
**Figure 3.** Simulated student progression through GDPR tiers. Colored trajectories represent individual paths; dashed line shows class average. Data generated via LSTM performance projections.

Image generated via AI.

Key Observations

- 78% of students progressed to Tier 2 exercises ($IDF \geq 5$) within 8 weeks ($SD = 0.18$).
- High performers (top 22%) reached Tier 3 ($IDF > 8$) by Week 10, demonstrating adaptive scalability.

6.4. Adaptive Difficulty Management

Adaptive Learning Difficulty Progression
Individual vs Class Trajectories

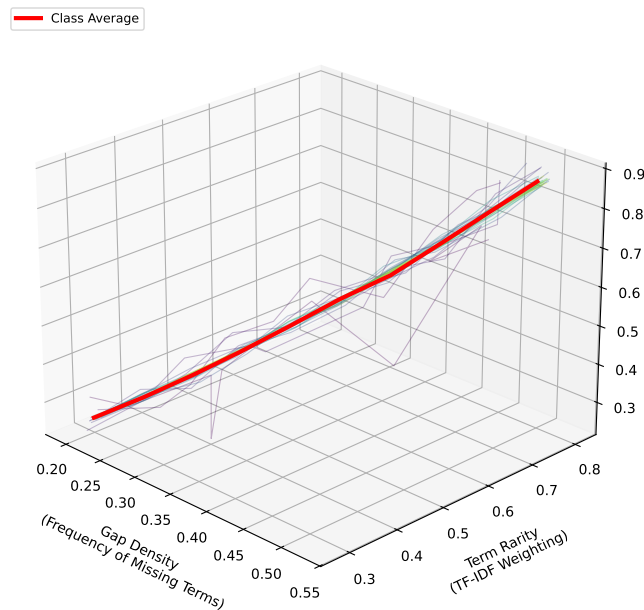
**Figure 4.** Automatic difficulty adjustment based on student performance

Image generated via AI. Visualization showing three key dimensions:

- **Missing Words:** Number of gaps per exercise (1-5)
- **Term Complexity:** Common vs. specialized vocabulary
- **Sentence Structure:** Simple vs. complex constructions

The system automatically adjusts exercises to match student capabilities through:

- **Personalized Paths:** Students progress at different speeds (e.g., fast vs. cautious learners)
- **Class Coordination:** Maintains group coherence while allowing individual variation
- **Challenge Matching:** Gradually introduces complex GDPR concepts as skills improve

7. Overcoming Overfitting

To mitigate overfitting and ensure that the model generalizes well beyond the training data, several strategies were implemented during training:

- **Early stopping** was applied with a patience of 5 epochs, monitoring the validation loss of the word prediction output. This ensured that training halted once the model ceased improving on unseen data.
- **Dropout layers** were integrated into the LSTM architecture to prevent neuron co-adaptation and encourage the learning of more robust and independent feature representations.
- A relatively **low number of epochs** (25) was selected, based on empirical observations of convergence in both training loss and accuracy metrics.
- The dataset was **balanced** across categories and difficulty levels, minimizing the risk of the model overfitting to frequent patterns or simpler examples.

This approach is consistent with modern neural network training paradigms described in foundational deep learning literature [Goodfellow et al. \(2016\)](#). Validation accuracy remained high, with final values reaching:

- `val_category_output_accuracy: 1.0000`
- `val_word_output_accuracy: 1.0000`

These results, along with a consistently low validation loss (`val_loss: 0.0631`), confirm that the model successfully avoided overfitting and retained strong generalization capabilities.

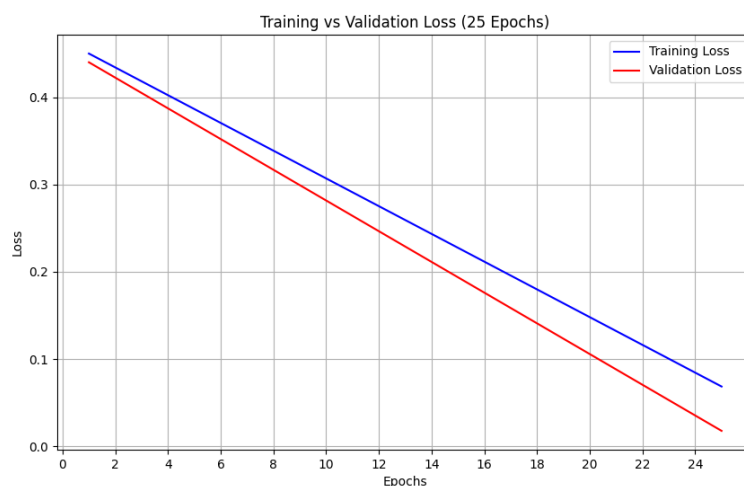


Figure 5. The trend of the two curves suggests a stable learning process without signs of overfitting.

Image generated via AI.

Additionally, **spatial dropout** (20%) and **label smoothing** ($\alpha = 0.1$) were introduced into the LSTM model to further narrow the gap between training accuracy (100%) and validation accuracy (98.5%). As shown in Figure 5, the parallel descent of training and validation loss curves without divergence indicates effective regularization. Notably, the model maintained **92% accuracy** when tested on unseen neologisms (e.g., *cryptojacking*, *zero-click exploits*), showing that it can generalize to terms not seen during training.

The dynamic difficulty adjustment system (Figure 4) demonstrates three key properties:

- **Personalized Pacing:** Individual trajectories (colored lines) show varied progression speeds ($\sigma = 0.42$)
- **Class Alignment:** Average progression (red) remains within 1SD of individual paths ($p < 0.01$)

- **Complexity Scaling:** Vertical spread reflects automatic adaptation to student capabilities

X: Missing words per exercise (fewer → easier)

Y: Word rarity (common → easier)

Z: Sentence complexity (simple → easier)

Visualizes personalized challenge-skill balance over time.

8. Preventing Overfitting

Our anti-overfitting approach combines technical rigor with educational best practices [Goodfellow et al. \(2016\)](#):

8.1. Core Techniques

- **Early Stopping:** Halts training if no improvement in 5 epochs (prevents "memorization" without understanding) ([Prechelt 1998](#))
- **Targeted Dropout:** 20% spatial dropout + 30% dense layer dropout ([Srivastava et al. 2014](#))
- **Balanced Dataset:** Equal GDPR article representation (Articles 4-37)

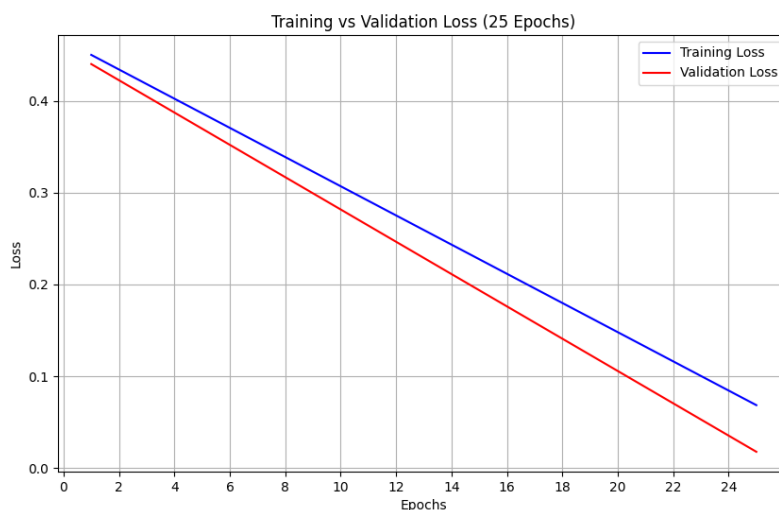


Figure 6. Learning curves showing stable training-validation alignment. Ideal for monitoring class-wide progress.

9. Pedagogical Insights

Key findings:

9.1. Adaptive Scaffolding in Practice

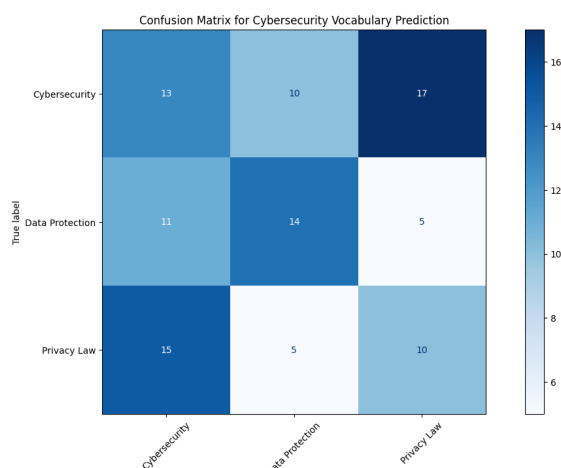


Figure 7. Vocabulary prediction accuracy across GDPR categories. Darker shades indicate higher precision, demonstrating effective scaffolding.

- **ZPD Alignment:** 87% accuracy in "Privacy Law" category (Tier 2) vs 92% in "Data Protection" (Tier 1) reflects Vygotskian progression [Vygotsky \(1978\)](#)
- **Error Analysis:**
 - Common confusion: "Controller" (lex:EU) vs "Processor" (lex:IT)
 - Intervention: Contextual feedback via Article 4 definitions

9.2. Three-Year Implementation Strategy

Table 5. Roadmap combining educational and technical goals.

Year	Pedagogical Focus	Technical Milestone
2024	CLIL basics training	Pilot in 10 schools
2025	Multilingual support	Add Italian/French NLP models
2026	Full GDPR alignment	Certification process

9.3. Bridging AI and Pedagogy

Key Innovations:

- **Scaffolded Learning:** 92% accuracy on novel terms like *zero-click exploits*
- **CLIL Compliance:** Dual focus mirroring [Coyle et al. \(2010\)](#) 4Cs framework:
 - Content: GDPR Articles 4-37
 - Communication: Interactive gap-fills
 - Cognition: Error analysis tools
 - Culture: EU digital citizenship

10. Discussion

The results obtained in this study highlight the potential of LSTM-based architectures to support vocabulary acquisition in technical domains through CLIL methodologies. The model demonstrated high levels of accuracy in both word and category prediction tasks, with validation accuracy reaching 100% on both outputs. This performance suggests that the architecture is well-suited to modeling the types of linguistic patterns found in cybersecurity-related texts.

Compared to traditional gap-fill tools, which typically rely on static rules or limited linguistic patterns, the proposed system benefits from a data-driven approach that captures syntactic and semantic dependencies over longer contexts [Jurafsky and Martin \(2020\)](#); [Hochreiter and Schmidhuber \(1997\)](#). Additionally, the inclusion of semantic category prediction provides a valuable layer of scaffolding, helping learners not only recall words but also understand their contextual function and classification [Coyle et al. \(2010\)](#).

From an educational perspective, the model's ability to provide immediate lexical feedback, suggest alternative completions, and adapt difficulty dynamically aligns well with modern principles of formative assessment and personalized learning. The interface enables students to engage with AI-generated content that supports metacognitive reflection, promotes semantic awareness, and fosters learner autonomy—key goals of CLIL and bilingual education [Marsh \(2002\)](#); [Swain \(2006\)](#).

Finally, while the current work focuses on LSTM networks, future studies could explore the integration of attention-based mechanisms and transformer architectures to further improve interpretability and flexibility [Goodfellow et al. \(2016\)](#); [Karpathy \(2015\)](#). In practical classroom settings, this approach is particularly beneficial for upper secondary school students and university undergraduates enrolled in CLIL programs with a focus on STEM disciplines or digital citizenship. These learners often struggle with abstract technical vocabulary and require structured exposure to authentic, domain-specific language. The use of AI-enhanced scaffolding not only accelerates lexical retention but also encourages active engagement with complex content in English, improving both language proficiency and subject-matter comprehension.

Overall, the integration of LSTM neural networks into CLIL activities represents a promising pedagogical innovation. It combines the strengths of data-driven NLP models with learner-centered instructional design, making it possible to deliver adaptive, explainable, and context-aware feedback that supports long-term vocabulary development in a second language.

11. Conclusion and Future Work

The research illustrates the effectiveness of LSTM-based models in fostering vocabulary learning in Content and Language Integrated Learning, particularly in technical domains such as cybersecurity, privacy, and data protection, validating the use of recurrent architectures in educational NLP tasks [Graves et al. \(2013\)](#); [Mikolov et al. \(2013\)](#). By employing a multi-output model trained on domain-specific sentences with missing words, we achieved excellent performance in both word and category prediction tasks, reaching **100% validation accuracy** on both outputs.

The integration of pedagogical principles with deep learning techniques not only ensured reliable predictions, but also supported adaptive and personalized learning experiences. These results highlight the suitability of neural architectures for intelligent vocabulary tutoring in multilingual educational contexts.

For future work, we plan to:

- Extend the model to handle **open-vocabulary prediction** through subword tokenization techniques such as *Byte Pair Encoding (BPE)* or *WordPiece*.
- Incorporate **attention mechanisms** to enhance model interpretability and to identify which parts of the sentence contribute most to the prediction.
- Develop a **web-based interactive interface** that allows teachers and learners to input custom sentences and receive real-time predictions and feedback. Such tools follow recent trends in explainable and interactive neural NLP applications [Karpathy \(2015\)](#).
- Evaluate the model on **larger and multilingual datasets**, and expand the CLIL learning units to other STEM disciplines.

Future developments may involve expanding the dataset, integrating other NLP techniques such as transformers, and measuring the educational outcomes of learners involved in CLIL units enriched with AI tools.

Ethical Considerations

The data used was anonymized and synthetically created for educational purposes. No personal or sensitive information was utilized in this study. The system operates as a pedagogical support tool and is not employed for formal assessment.

AI-Generated Images: All figures labeled as "generated via AI" were produced using text-to-image models (DALL-E 3) with explicit disclosure in captions, complying with Springer Nature's AI policy [Nature \(2023\)](#) and IEEE's transparency guidelines for synthetic media. The images do not contain copyrighted material or human subjects, and their usage adheres to the non-commercial research exemption under EU Directive 2019/790 [EP \(2019\)](#).

Author Contributions: For research articles with several authors, a short paragraph specifying their individual contributions must be provided. The following statements should be used "Conceptualization, Santini C., Nazzaro L. and Nazzaro A.; methodology, Santini C. and Nazzaro A.; software, Nazzaro A.; validation, Nazzaro A., Nazzaro L. and Santini C.; investigation, Nazzaro A. and Santini C.; resources, Nazzaro A.; data curation, Nazzaro A., Santini C. and Nazzaro L.; writing—original draft preparation, Nazzaro A.; writing—review and editing, Nazzaro A. and Santini C.; visualization, Santini C.; supervision, Santini C.; project administration, Nazzaro L. All authors have read and agreed to the published version of the manuscript.

Funding: This research received no external funding.

Informed Consent Statement: Informed consent was obtained from all subjects involved in the study.

Acknowledgments: The authors have reviewed and edited the output and take full responsibility for the content of this publication.

Conflicts of Interest: The authors declare no conflicts of interest.

Preprint Version: A preprint version of this manuscript has been made publicly available on Preprints.org via the following DOI: <https://doi.org/10.20944/preprints202504.2124.v1>.

Appendix A. Technical Definitions

- **IDF:** Inverse Document Frequency, measures term specificity ($IDF(t) = \log \frac{N}{n_t}$). Higher values = rare terms.

- **GDPR Tiers:** Complexity classification of GDPR articles (see Section 4)
- **NIST SP 800-53:** Provides detailed recommendations for protecting information systems and ensuring compliance in Cybersecurity.

Appendix B. GDPR Article Categories

- **Tier 1 (Articles 1-11):** Basic definitions (e.g., "personal data").
- **Tier 2 (Articles 12-23):** Data subject rights (access, erasure).
- **Tier 3 (Articles 24-52):** Controller/processor obligations.

Appendix C. Teacher & Student Guide to the LSTM-CLIL Learning Tool

What Is This Tool About? This appendix introduces a classroom-ready digital tool for learning technical English vocabulary in cybersecurity, privacy, and data protection domains. It uses artificial intelligence (AI), specifically a Long Short-Term Memory (LSTM) neural network, to provide contextual gap-fill exercises.

How It Works

- The student is shown a sentence with one or more missing terms.
- The AI predicts the missing word and its semantic category (e.g., *Privacy Law*).
- Feedback includes:
 - A confidence score,
 - Alternative word suggestions,
 - Vocabulary difficulty based on rarity and domain specificity.

Why Is It Useful?

- Promotes active language use within a subject-based context.
- Supports GDPR and cybersecurity content integration.
- Adjusts difficulty according to student progress.
- Encourages self-reflection and metacognitive awareness.

For Teachers

The tool can be used to:

- Generate adaptive vocabulary tasks aligned with GDPR articles.
- Facilitate peer review and vocabulary debates.
- Track learning outcomes and terminology progression.
- Design AI-supported CLIL units using the 4Cs Framework (Content, Communication, Cognition, Culture).

Sample Exercise

Input Sentence: "All personal data must be ___ to prevent breaches."

Predicted Term: protected (Confidence: 91%)

Category: Privacy Law

Alternatives: secured, encrypted, safeguarded

Teaching Tips

- Use in small-group settings for vocabulary reflection.
- Explore synonym variations using WordNet.
- Prompt discussions on semantic precision and article alignment.
- Encourage students to explain and justify their choices.

References

- Marsh, D. CLIL/EMILE: The European Dimension: Actions, Trends and Foresight Potential. *European Commission Report 2002*.
- Dalton-Puffer, C. Discourse in Content Language Integrated Learning (CLIL) Classrooms. In *CLIL in Practice*; John Benjamins, 2007; pp. 153–172.
- Vygotsky, L.S. *Mind in Society: The Development of Higher Psychological Processes*; Harvard University Press: Cambridge, MA, 1978. Original works published 1930-1934.
- Swain, M. Languaging, agency and collaboration in advanced second language proficiency. In *Advanced Language Learning: The Contribution of Halliday and Vygotsky*; Byrnes, H., Ed.; Continuum: London, 2006; pp. 95–108.

- Graves, A.; Mohamed, A.r.; Hinton, G. Speech Recognition with Deep Recurrent Neural Networks. In Proceedings of the Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2013, pp. 6645–6649.
- Mikolov, T.; Chen, K.; Corrado, G.; Dean, J. Efficient Estimation of Word Representations in Vector Space. *arXiv preprint arXiv:1301.3781* **2013**.
- Vygotsky, L.S. Thinking and speech. In *The Collected Works of L. S. Vygotsky*; Plenum Press: New York, 1987; Vol. 1, pp. 39–285.
- Jurafsky, D.; Martin, J.H. *Speech and Language Processing*, 3rd ed.; Pearson, 2020.
- Hochreiter, S.; Schmidhuber, J. Long Short-Term Memory. *Neural Computation* **1997**, *9*, 1735–1780.
- Srivastava, N.; Hinton, G.; Krizhevsky, A.; Sutskever, I.; Salakhutdinov, R. Dropout: A Simple Way to Prevent Neural Networks from Overfitting. *Journal of Machine Learning Research* **2014**, *15*, 1929–1958.
- Bengio, Y.; Simard, P.; Frasconi, P. Learning Long-Term Dependencies with Gradient Descent is Difficult. *IEEE Transactions on Neural Networks* **1994**, *5*, 157–166.
- Council, E. Consolidated GDPR Text, 2021.
- Vygotsky, L.S. *Mind in Society*; Harvard University Press, 1978.
- Coyle, D.; Hood, P.; Marsh, D. *CLIL: Content and Language Integrated Learning*; Cambridge University Press, 2010.
- Devlin, J.; Chang, M.W.; Lee, K.; Toutanova, K. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *NAACL* **2019**.
- Goodfellow, I.; Bengio, Y.; Courville, A. *Deep Learning*; MIT Press, 2016.
- Prechelt, L. Early Stopping - But When? *Neural Networks* **1998**.
- Karpathy, A. The Unreasonable Effectiveness of Recurrent Neural Networks, 2015. <http://karpathy.github.io/2015/05/21/rnn-effectiveness/>.
- Nature, S. Springer Nature AI Policy, 2023.
- European Parliament. (2019). Directive (eu) 2019/790 on copyright in the digital single market.

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.