

# <sup>1</sup> LAiSER: A Taxonomy-Aware Framework for Skill Extraction and Research

<sup>3</sup> **Satya Phanindra Kumar Kalaga**  <sup>1</sup> and **Bharat Khandelwal**  <sup>1</sup>

<sup>4</sup> <sup>1</sup> Program on Skills, Credentials and Workforce Policy, Institute of Public Policy, The George  
<sup>5</sup> Washington University, USA

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

## Software

- <sup>6</sup> [Review](#) 
- <sup>7</sup> [Repository](#) 
- <sup>8</sup> [Archive](#) 

Editor: 

Submitted: 01 November 2025

Published: unpublished

## License

Authors of papers retain copyright<sup>15</sup> and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](https://creativecommons.org/licenses/by/4.0/)).<sup>16</sup>

## <sup>6</sup> Summary

<sup>7</sup> LAiSER is an artificial intelligence framework for extracting, standardizing, and aligning skill information from unstructured text with established skill taxonomies. The system addresses <sup>8</sup> the lack of standardized, machine-readable skill information needed to facilitate communication <sup>9</sup> between learners, educators, and employers. LAiSER employs a two-stage pipeline <sup>10</sup> combining large language models (LLMs) with semantic vector search using FAISS (Johnson <sup>11</sup> et al., 2019) for high-precision skill extraction and taxonomy alignment to standards <sup>12</sup> like ESCO (European Commission, 2020). The framework supports multiple computational <sup>13</sup> backends—GPU-accelerated inference, cloud APIs, and CPU-only environments—making it <sup>14</sup> accessible across diverse research and operational contexts.

## <sup>14</sup> Statement of need

<sup>15</sup> The contemporary labor market is characterized by rapid technological change, evolving <sup>16</sup> skill requirements, and growing disconnect between educational curricula and industry needs <sup>17</sup> (Autor & Dorn, 2013). Traditional skill extraction approaches rely on manual annotation, <sup>18</sup> expert judgment, or keyword matching—methods that are labor-intensive, error-prone, and <sup>19</sup> fail to capture semantic richness (Boselli et al., 2018). LAiSER fills this gap by providing <sup>20</sup> an accessible, flexible, and accurate framework for automated skill extraction and taxonomy <sup>21</sup> alignment, enabling researchers, educators, and workforce professionals to analyze skill demands <sup>22</sup> at scale.

## <sup>25</sup> Key Features

<sup>26</sup> LAiSER provides three core innovations:

- <sup>27</sup> **1. Intelligent Text Preprocessing:** Domain-aware preprocessing removes extraneous information (company branding, legal boilerplate, benefits) while preserving task-relevant skill content.
- <sup>28</sup> **2. Multi-Model Skill Extraction:** Flexible architecture supporting multiple LLMs (vLLM, HuggingFace Transformers, Gemini API) with automatic fallback mechanisms for deployment across different computational environments.
- <sup>29</sup> **3. Taxonomy-Aware Alignment:** Semantic similarity-based alignment using sentence transformers (Reimers & Gurevych, 2019) and FAISS vector search maps extracted skills to <sup>30</sup> standardized taxonomies, enabling cross-domain skill analysis.

<sup>31</sup> The framework also supports extraction within the Knowledge, Skills, and Abilities (KSA) framework <sup>32</sup> with automatic proficiency level classification using the Scottish Credit and Qualifications <sup>33</sup>

<sup>38</sup> Framework (SCQF) ([Scottish Credit and Qualifications Framework Partnership, 2019](#)).

## <sup>39</sup> Implementation

<sup>40</sup> LAiSER is implemented in Python 3.9+ using PyTorch, HuggingFace Transformers, FAISS,  
<sup>41</sup> spaCy, Sentence-Transformers, and pandas. The modular service architecture separates data  
<sup>42</sup> access, skill extraction, and taxonomy alignment into loosely coupled layers.

<sup>43</sup> The system accepts pandas DataFrames containing textual descriptions and produces structured  
<sup>44</sup> output including raw extracted skills, taxonomy-aligned canonical skills, taxonomy identifiers  
<sup>45</sup> (e.g., ESCO codes), and semantic similarity scores. For KSA extraction, output includes SCQF  
<sup>46</sup> proficiency levels, knowledge requirements, and task abilities.

## <sup>47</sup> Applications

<sup>48</sup> LAiSER enables diverse applications: labor market intelligence through large-scale job adver-  
<sup>49</sup> tisement analysis; curriculum development through skill gap identification between educational  
<sup>50</sup> programs and industry demands; skills-based hiring through standardized job-candidate match-  
<sup>51</sup> ing; career pathway analysis through skill similarity assessment across occupations; and  
<sup>52</sup> workforce policy research through regional skills assessments.

## <sup>53</sup> Availability

<sup>54</sup> The source code is available at <https://github.com/LAiSER-Software/extract-module> under  
<sup>55</sup> the MIT License. Installation via pip:

```
pip install laiser[gpu] # GPU-accelerated
pip install laiser[cpu] # CPU-only
```

<sup>56</sup> Example usage:

```
from laiser.skill_extractor_refactored import SkillExtractorRefactored
import pandas as pd

extractor = SkillExtractorRefactored(model_id="gemini", api_key="YOUR_API_KEY", use_gpu=
data = pd.read_csv("job_descriptions.csv")
results = extractor.extract_and_align(data, id_column="job_id",
text_columns=["description"], input_type="job_desc")
```

## <sup>57</sup> Acknowledgments

<sup>58</sup> The authors acknowledge the George Washington University Institute of Public Policy and the  
<sup>59</sup> Program on Skills, Credentials, and Workforce Policy for institutional support. This project  
<sup>60</sup> was supported by grants from the Walmart Foundation and the Gates Foundation. The authors  
<sup>61</sup> thank the GW Open Source Program Office and the developers of HuggingFace Transformers,  
<sup>62</sup> FAISS, and spaCy.

## <sup>63</sup> References

<sup>64</sup> Autor, D. H., & Dorn, D. (2013). The growth of low-skill service jobs and the polarization of  
<sup>65</sup> the US labor market. *American Economic Review*, 103(5), 1553–1597. <https://doi.org/10.1257/aer.103.5.1553>

- 67 Boselli, R., Cesarini, M., Mercorio, F., & Mezzanica, M. (2018). Classifying online job  
68 advertisements through machine learning. *Future Generation Computer Systems*, 86,  
69 319–328. <https://doi.org/10.1016/j.future.2018.03.035>
- 70 European Commission. (2020). *ESCO: European skills, competences, qualifications and  
71 occupations*. <https://ec.europa.eu/esco>
- 72 Johnson, J., Douze, M., & Jégou, H. (2019). Billion-scale similarity search with GPUs.  
73 *IEEE Transactions on Big Data*, 7(3), 535–547. <https://doi.org/10.1109/TB DATA.2019.2921572>
- 75 Reimers, N., & Gurevych, I. (2019). Sentence-BERT: Sentence embeddings using siamese  
76 BERT-networks. *Proceedings of the 2019 Conference on Empirical Methods in Natural  
77 Language Processing and the 9th International Joint Conference on Natural Language  
78 Processing (EMNLP-IJCNLP)*, 3982–3992. <https://doi.org/10.18653/v1/D19-1410>
- 79 Scottish Credit and Qualifications Framework Partnership. (2019). *SCQF handbook: User  
80 guide*. <https://scqf.org.uk>