

# So sánh các bộ dữ liệu Fact-checking và Suy luận ngôn ngữ tự nhiên tiếng Việt

## Bảng so sánh chi tiết các bộ dữ liệu

Bảng dưới đây tóm tắt đặc trưng của các bộ dữ liệu: **ViAdverNLI** (bao gồm 3 vòng R1, R2, R3), **ViNLI**, **ViWikiFC**, **ViFactCheck** và **ISE-DSC01**, dựa trên các tiêu chí đã liệt kê:

Tiêu chí	ViAdverNLI (R1-R3)	ViNLI	ViWikiFC	ViFactCheck	ISE-DSC01
Nguồn dữ liệu	<p>Các <b>câu premise</b> trích từ báo điện tử VnExpress (700+ tin tức đa chủ đề) <sup>1</sup> <sup>2</sup> .</p> <p><b>Hypothesis</b> do người chú thích viết để đánh lừa mô hình NLI.</p>	<p>Câu <b>premise</b> lấy từ &gt;800 bài báo trực tuyến (13 chủ đề: thời sự, giải trí, thể thao, v.v.) <sup>3</sup> .</p> <p><b>Hypothesis</b> do người tạo theo hướng dẫn NLI tiêu chuẩn.</p>	<p><b>Câu bằng chứng</b> (evidence) lấy từ <b>Wikipedia tiếng Việt</b> – trích các phát biểu thực tế từ bài wiki rồi chuyển thành <i>claim</i> cần kiểm chứng <sup>4</sup> .</p>	<p><b>Claim</b> (phát biểu) và <b>đoạn evidence</b> (bằng chứng) trích từ <b>bài báo tin tức Việt Nam</b> uy tín (Thanh Niên, Tuổi Trẻ, VnExpress, ...) thuộc 12 lĩnh vực khác nhau <sup>5</sup> . Mỗi <i>claim</i> gắn với bối cảnh bài báo thực tế.</p>	<p><b>Claim</b> và <b>đoạn văn ngữ cảnh</b> lấy từ <b>nhiều bài báo tiếng Việt</b> đa lĩnh vực (thời sự, kinh tế, khoa học, v.v.) <sup>6</sup> <sup>7</sup> . Nguồn tin chọn từ báo phổ biến, đáng tin cậy, đảm bảo tính đa dạng và thực tế.</p>

Tiêu chí	ViAdverNLI (R1–R3)	ViNLI	ViWikiFC	ViFactCheck	ISE-DSC01
<b>Số lượng mẫu</b>  (train/ dev/test)	~ <b>10.0k</b> cặp premise– hypothesis (tổng hợp 3 vòng) <sup>2</sup> . Chia 8:1:1 thành 8,012 train / 1,000 dev / 1,000 test <sup>8</sup> . (Mỗi vòng ~3.3k cặp: R1 ~2.6k train, 330 dev, 330 test; R2 ~2.7k train, 330 dev, 330 test; R3 ~2.7k train, 340 dev, 340 test) <sup>9</sup> .	> <b>30k</b> cặp premise– hypothesis có nhãn (khoảng 24k train, 3k dev, 3k test – chia 80/10/10%) <sup>10</sup> <sup>11</sup> . Mỗi premise đi kèm 8 hypothesis (2 cho mỗi nhãn) nhằm cân bằng dữ liệu <sup>12</sup> .	> <b>20k</b> claim được gán nhãn kèm bằng chứng (≈16k train / 2k dev / 2k test) <sup>13</sup> . Phân bố nhãn <b>cân bằng</b> giữa 3 lớp <sup>14</sup> . Mỗi claim có một hoặc vài câu bằng chứng tương ứng từ wiki.	<b>7,232</b> cặp claim– evidence (5,126 train / 1,026 dev / 1,080 test – ~70/10/20%) <sup>15</sup> .	~ <b>49.7k</b> cặp claim–đoạn văn (38,684 train / 4,793 dev / 5,396 test) <sup>7</sup> – một trong những bộ dữ liệu fact- checking lớn nhất cho tiếng Việt. Nhãn phân bố gần như <b>cân bằng</b> giữa các lớp <sup>16</sup> .

Tiêu chí	ViAdverNLI (R1-R3)	ViNLI	ViWikiFC	ViFactCheck	ISE-DSC01
Định dạng dữ liệu	Mỗi mẫu gồm: một <b>premise</b> (câu gốc từ tin tức) + một <b>hypothesis</b> (câu do annotator viết) + một <b>nhãn suy luận</b> (Entail/Contr/Neutral). Được xây dựng qua nhiều vòng điều chỉnh để tăng độ khó cho mô hình <sup>17</sup> <sup>18</sup> .	Mỗi mẫu gồm: một câu <b>premise</b> và một câu <b>hypothesis</b> , kèm <b>nhãn suy luận logic</b> giữa chúng (kéo theo, mâu thuẫn hoặc trung tính) <sup>19</sup> <sup>20</sup> . Định dạng tuân theo chuẩn NLI (giống SNLI/MultiNLI).	Mỗi mẫu gồm: một <b>claim</b> (câu phát biểu cần kiểm chứng) + <b>bộ bằng chứng</b> (1 hoặc nhiều câu Wikipedia liên quan) + <b>nhãn</b> đánh giá độ đúng/sai của claim dựa trên bằng chứng <sup>21</sup> <sup>22</sup> . Định dạng tương tự FEVER cho tiếng Anh.	Mỗi mẫu gồm: một <b>claim</b> và một <b>đoạn evidence</b> (trích từ bài báo liên quan) + <b>nhãn</b> xác định tính đúng/sai của claim dựa trên evidence đó <sup>23</sup> <sup>24</sup> . Định dạng tương tự bài toán kiểm chứng claim FEVER nhưng áp dụng cho ngữ liệu báo chí Việt <sup>24</sup> .	Mỗi mẫu gồm: một <b>claim</b> và một <b>đoạn văn ngữ cảnh</b> liên quan + <b>nhãn</b> xác định claim được hỗ trợ, bác bỏ hay không đủ thông tin từ đoạn văn <sup>25</sup> . (Trong tập huấn luyện, đoạn văn thường chứa bằng chứng trực tiếp/gián tiếp; trong thi đấu, mô hình phải tự tìm câu chứng cứ từ văn bản dài) <sup>26</sup> <sup>27</sup> .

Tiêu chí	ViAdverNLI (R1-R3)	ViNLI	ViWikiFC	ViFactCheck	ISE-DSC01
Bộ nhãn & phân bố	<p><b>3 nhãn NLI:</b> Entailment (kéo theo), Contradiction (mâu thuẫn), Neutral (trung tính) – tương ứng Hỗ trợ, Phản bác, Không đủ thông tin (theo ngữ cảnh) <sup>28</sup> .</p> <p><b>Phân bố:</b> gần <b>cân bằng</b> tổng thể (mỗi nhãn ~33%) nhưng thay đổi theo vòng: ví dụ vòng 1 ít NEI hơn (chủ yếu Entail/Contr) so với vòng 3 (NEI tăng) – phản ánh chiến lược tạo câu ngày càng đa dạng.</p>	<p><b>3 nhãn NLI tiêu chuẩn:</b> Entailment (kéo theo), Contradiction (mâu thuẫn), Neutral (trung tính) <sup>20</sup> . Nhờ thiết kế mỗi premise sinh 2 hypothesis mỗi loại, tập dữ liệu được cân bằng 3 nhãn (xấp xỉ 1/3 mỗi lớp).</p>	<p><b>3 nhãn:</b> <i>Supports</i> (có căn cứ), <i>Refutes</i> (bác bỏ), <i>Not Enough Information</i> (thiếu thông tin) – theo schema FEVER <sup>29</sup> <sup>30</sup> . Phân bố nhãn đã được <b>cân bằng</b> giữa các lớp (~33% mỗi loại) trên toàn bộ tập <sup>14</sup> .</p>	<p><b>3 nhãn:</b> <i>Support</i> (Hỗ trợ), <i>Refute</i> (Bác bỏ), <i>Not Enough Information</i> (NEI) (Không đủ thông tin) <sup>31</sup> <sup>32</sup> . Mỗi bài báo gốc tạo đúng 2 claim mỗi nhãn, nên tập dữ liệu <b>cân bằng hoàn hảo</b> (mỗi lớp ~1/3 tổng số claim) <sup>33</sup> <sup>34</sup> . (Kappa đồng thuận giữa các annotator = 0,83 – rất cao) <sup>35</sup> .</p>	<p><b>3 nhãn:</b> <i>Supported</i> (Được hỗ trợ), <i>Refuted</i> (Bị bác bỏ), <i>NEI</i> (Không đủ thông tin) <sup>36</sup> <sup>37</sup> .</p> <p>Nhãn được xác định qua đối chiếu claim với nội dung bài báo: nếu bài báo khẳng định claim ⇒ <i>Supported</i>; mâu thuẫn claim ⇒ <i>Refuted</i>; không đề cập đủ ⇒ <i>NEI</i> <sup>38</sup> <sup>39</sup> .</p> <p>Phân bố nhãn ~cân bằng (vd. train: 12.8k Sup, 12.6k Ref, 13.3k NEI) <sup>40</sup> .</p>

Tiêu chí	ViAdverNLI (R1-R3)	ViNLI	ViWikiFC	ViFactCheck	ISE-DSC01
Độ dài TB (claim & evidence)	<p><b>Premise (tin tức):</b> ~24 từ (tương tự ViNLI) <sup>41</sup>.</p> <p><b>Hypothesis:</b> ngắn, trung bình ~12-15 từ (ngắn hơn hẳn so với ViNLI ~18 từ) <sup>42</sup>. Các câu hypothesis thường súc tích, tập trung thay đổi nhỏ để đánh lừa mô hình.</p>	<p><b>Premise:</b> ~24,5 từ;</p> <p><b>Hypothesis:</b> ~18,1 từ (trung bình) <sup>11</sup>.</p> <p>Hypothesis thường ngắn hơn premise, nhưng vẫn đủ rõ nghĩa cho suy luận <sup>43</sup>. (Ngắn nhất 4 từ, dài nhất 68 từ; đa số 10-23 từ) <sup>44</sup>.</p>	<p><b>Claim:</b> thường là một câu ngắn gọn (~15-20 từ). <b>Bằng chứng:</b> mỗi câu bằng chứng Wikipedia ~20 từ; một claim có thể kèm 1-2 câu bằng chứng nên tổng độ dài ngữ cảnh ~20-40 từ. (Suy ra từ độ dài câu wiki trung bình, tương tự FEVER).</p>	<p><b>Claim:</b> ngắn, thường 1 câu (~12-15 từ).</p> <p><b>Evidence:</b> đoạn trích báo ~1-3 câu liên quan (có thể ~30-50 từ). Do mỗi claim chỉ dựa vào một đoạn chứng cứ cụ thể nên ngữ cảnh khá gọn. (Theo mô tả, evidence là "đoạn trích nội dung báo liên quan", thường chỉ vài dòng) <sup>23</sup>.</p>	<p><b>Claim:</b> thường 1 câu (~10-20 từ). <b>Đoạn ngữ cảnh:</b> trung bình dài hơn evidence của ViFactCheck, có thể là một đoạn văn (~50-100 từ). Trong dataset train, đoạn văn đã được cắt để chứa thông tin chính yếu (tránh quá dài) <sup>26</sup>. Trong chế độ thi, mô hình phải xử lý cả bài báo dài nên độ dài ngữ cảnh thực tế có thể lớn.</p>

## Quy mô từ vựng & đa dạng ngôn ngữ

Tương đối **đa dạng, nhiều từ mới**:

Annotator cố ý **không lặp từ** của premise trong hypothesis – tỷ lệ trùng từ thấp và thường dùng từ đồng nghĩa, cách diễn đạt khác <sup>45</sup>. Mỗi vòng bổ sung từ vựng mới (đặc biệt là danh từ, động từ) để mô hình khó suy luận dựa trên từ khóa bề mặt <sup>46</sup>. Chủ đề tin tức phong phú (thời sự, giải trí, khoa học,...).

Quy mô **lớn, phủ nhiều chủ đề** (13 lĩnh vực tin tức) nên từ vựng khá phong phú <sup>3</sup>. Hướng dẫn sinh dữ liệu yêu cầu **không sao chép y nguyên**

premise mà dùng từ ngữ của riêng mình <sup>48</sup> – do đó corpus chứa nhiều cách diễn đạt lại, bao gồm từ đồng nghĩa, câu chủ-bị động, thay đổi trạng từ, v.v. (theo các luật tạo dữ liệu) <sup>49</sup>.

<sup>50</sup>.

**Đa dạng kiến thức bách khoa**:

bao gồm nhiều tên riêng, thuật ngữ về địa danh, nhân vật, sinh vật, lịch sử... do lấy từ Wikipedia. Văn phong bách khoa chung, nhưng nhờ cơ chế tạo claim (đối chi tiết hoặc thêm thông tin ngoài phạm vi) nên có xuất hiện cả câu đúng lẫn sai, thông tin hư cấu. Từ vựng bao trùm các lĩnh vực có trên Wikipedia tiếng Việt (khá rộng, nhưng bị giới hạn ở phạm vi tri thức wiki).

**Đa miền ngôn ngữ**

**báo chí**: bao quát 12 lĩnh vực từ chính trị, y tế đến giải trí <sup>5</sup> nên tập hợp từ vựng rất **phong phú**. Các claim và evidence chứa nhiều số liệu, tên riêng, thuật ngữ chuyên ngành (phù hợp từng lĩnh vực). Tuy quy mô mẫu vừa phải (~7k) nhưng do tạo 6 claim mỗi bài báo theo kịch bản khác nhau <sup>33</sup>, ngôn ngữ biểu đạt khá đa dạng (bao gồm diễn đạt lại sự kiện thực tế, cố ý sửa chi tiết tạo thông tin sai, và đặt câu hỏi ngoài phạm vi bài).

**Rất lớn và phong phú**: gần 50k mẫu từ hàng nghìn bài báo trải rộng nhiều chuyên mục, do đó **quy mô từ vựng lớn nhất** trong các bộ so sánh. Văn bản ngữ cảnh dài hơn, chứa câu phức và thông tin nền. Đa số từ vựng là ngôn ngữ báo chí phổ thông, nhưng nhờ lượng dữ liệu lớn, mô hình học được nhiều cách diễn đạt khác nhau của cùng một nội dung. Tập dữ liệu cũng bao gồm nhiều **tên riêng, thuật ngữ chuyên ngành** từ các mảng khác nhau (kinh tế, khoa học, thể thao,...).

## Tính chất *adversarial*

**Có** – Bộ dữ liệu *adversarial* duy nhất: Được xây dựng theo quy trình **human-and-model-in-the-loop** qua 3 vòng <sup>17</sup>. Mỗi vòng, annotator tìm cách viết hypothesis gây **hiểu lầm cho mô hình** NLI hiện tại, sau đó lọc những mẫu mô hình dự đoán sai (xác nhận bởi người) làm dữ liệu vòng tiếp theo <sup>51</sup> <sup>52</sup>. Qua từng vòng, mô hình ngày càng mạnh nhưng vẫn bị đánh lừa bởi câu khó hơn. Kết quả, ViAdverNLI chứa nhiều mẫu hóc búa mà **mô hình SOTA chỉ đạt ~48% độ chính xác** trên test <sup>53</sup> – thấp hơn hẳn các dataset khác.

**Không (chuẩn)** – Được xây dựng theo quy trình truyền thống, chú trọng chất lượng hơn là đánh lừa mô hình. Các cặp câu được viết và gán nhãn cẩn thận, tránh cặp quá đơn giản hay mơ hồ <sup>54</sup> <sup>55</sup>, nhưng **không có yếu tố adversarial chủ đích**. Mục tiêu là phản ánh suy luận thực tế, không tập trung khai thác lỗ hổng của mô hình.

**Không (theo chuẩn FEVER)** – Claim được tạo thủ công từ câu wiki nhưng không dựa trên phản hồi của mô hình. Tuy nhiên claim đòi hỏi suy luận phức tạp (cần kết hợp kiến thức hoặc phát hiện thông tin bị đảo), bộ dữ liệu **không sử dụng mô hình trong vòng tạo dữ liệu**. Vì thế, tính adversarial chỉ ở mức thay đổi thông tin để đánh đố người/máy theo kiểu FEVER, chưa phải adversarial multi-round.

**Không hẳn (nhưng có yếu tố sáng tạo)** – Dù không có mô hình trong vòng tạo dữ liệu, nhóm tác giả **cố tình tạo claim sai bằng cách chỉnh sửa chi tiết** và claim NEI bằng cách đưa thông tin ngoài bài <sup>34</sup> <sup>56</sup>. Do đó, dataset có những mẫu “gây nhiễu” giống thật (adversarial đối với người đọc). Tuy nhiên, các claim này được thiết kế nhằm đảm bảo đủ thông tin để con người nhận biết (dựa vào bài báo), chứ không nhằm đánh bại mô hình cụ thể nào.

**Không (bán tự động)** – Dữ liệu được tạo ra bằng **đối chiếu tự động + kiểm duyệt thủ công**, không phải qua nhiều vòng tương tác người-máy. Mục tiêu chính là quy mô lớn và đủ độ khó tổng quan, không nhằm đến việc tìm lỗi mô hình cụ thể. Tuy nhiên, do có khâu tự động, có thể chứa một số trường hợp nhiễu tự nhiên (mô hình dễ nhầm) – nhưng đây là tác dụng phụ hơn là thiết kế có chủ đích.

## Phương pháp gán nhãn

### Kết hợp con người & mô hình:

Người viết hypothesis, mô hình dự đoán, sau đó người xác nhận và lọc. Cụ thể, vòng 1 dùng XLM-R gán nhãn tạm cho mẫu annotator viết; vòng 2-3 dùng InfoXLM (huấn luyện trên dữ liệu mở rộng) đánh giá mẫu mới <sup>1</sup> <sup>57</sup>. Chỉ giữ lại những cặp mà mô hình dự đoán *sai* nhưng được người xác minh là hợp lệ để đưa vào tập dữ liệu <sup>51</sup> <sup>52</sup>. Cuối cùng, toàn bộ ~10k cặp được gán nhãn vàng bởi con người (theo đa số phiếu nếu cần).

### Gán nhãn thủ công chất lượng

**cao:** Các cặp premise-hypothesis do người viết theo hướng dẫn chi tiết (có bộ quy tắc tạo câu cho entailment và contradiction) <sup>58</sup> <sup>49</sup>. Mọi mẫu (nhất là dev/test) được **5 annotator gán nhãn độc lập**, chọn nhãn vàng bằng bỏ phiếu đa số <sup>59</sup> <sup>60</sup>. Kết quả 99,4% cặp được  $\geq 3/5$  phiếu trùng, chất lượng rất cao <sup>61</sup>.

### Gán nhãn thủ công theo FEVER:

Annotator được huấn luyện viết claim dựa trên câu wiki gốc rồi tự gán nhãn cho claim đó luôn (dựa trên bằng chứng wiki) <sup>62</sup> <sup>63</sup>. Mỗi claim sau đó đều được kiểm tra chéo với câu nguồn để đảm bảo nhãn đúng và câu claim hợp lệ. Bộ dữ liệu được nhóm tác giả soát xét nhằm loại bỏ lỗi và cân bằng tập. (Hiện tại dữ liệu được công bố có thể qua liên hệ nhóm tác giả) <sup>64</sup>.

### Gán nhãn thủ công nghiêm ngặt, có giám sát chuyên gia:

Chọn các bài báo, mỗi bài giao cho annotator tạo **6 cặp claim-evidence** (2 Hỗ trợ, 2 Bác bỏ, 2 NEI) dựa trên nội dung bài <sup>33</sup> <sup>56</sup>. Cụ thể: claim *Support* lấy sự kiện **có thật** trong bài, viết lại thành câu đúng; claim *Refute* **sửa một chi tiết** trong bài (số liệu, tên, thời gian...) tạo câu sai; claim *NEI* nêu thông tin **liên quan chủ đề nhưng không có trong bài** <sup>34</sup> <sup>56</sup>. Mỗi claim được đối chiếu lại với bài và gán nhãn + chỉ rõ đoạn bằng chứng. Nhờ quy trình kỹ lưỡng và

### Bán tự động + hiệu chỉnh thủ công:

Với quy mô lớn, ban đầu dùng kỹ thuật **cross-check tự động**: đối chiếu claim với bài báo gốc, tự động gán nhãn Supported nếu bài chứa thông tin khẳng định claim, Refuted nếu mâu thuẫn, NEI nếu không đề cập đủ <sup>38</sup> <sup>39</sup>. Sau đó, một phần dữ liệu được **kiểm tra thủ công** để sửa lỗi và đảm bảo cân bằng nhãn <sup>65</sup> <sup>66</sup>. Kết quả cuối đạt bộ dữ liệu ~50k với nhãn phân bố gần đều (vd. train mỗi nhãn ~12-13k) <sup>40</sup>.



Tiêu chí	ViAdverNLI (R1-R3)	ViNLI	ViWikiFC	ViFactCheck	ISE-DSC01
				huấn luyện thống nhất, dữ liệu đạt chất lượng cao (Fleiss' Kappa = 0,83) <sup>35</sup> .	

## Mục tiêu ứng dụng

Cung cấp **benchmark NLI adversarial** đầu tiên cho tiếng Việt, nhằm thử thách mô hình và cải thiện độ **robust**. Các mô hình SOTA khi đánh giá trên ViAdverNLI cho kết quả rất thấp (mô hình mạnh nhất chỉ ~48% trên test) <sup>53</sup>, cho thấy những điểm yếu cần khắc phục. Ngoài ra, khi huấn luyện trên ViAdverNLI, mô hình **cải thiện đáng kể hiệu quả trên các dataset NLI khác** <sup>53</sup> – chứng tỏ tập dữ liệu này giúp nâng cao khả năng khái quát và chống “bẫy” của mô hình. ViAdverNLI đặt nền móng cho hướng nghiên cứu NLI chống lại mẫu nhiễu,

Thiết lập **benchmark NLI tiếng Việt đầu tiên** (COLING 2022) <sup>67</sup>, làm chuẩn đánh giá các mô hình suy luận tiếng Việt. ViNLI đã được sử dụng rộng rãi để so sánh hiệu quả mô hình và thúc đẩy nghiên cứu chuyên biệt (ví dụ, phát sinh ViHealthNLI cho y tế) <sup>68</sup>. Nhờ ViNLI, cộng đồng có điểm bắt đầu để phát triển các mô hình NLI cho tiếng Việt, khi trước đó hầu như chưa có dữ liệu suy luận công khai nào <sup>69</sup>.

Cung cấp **benchmark fact-checking trên tri thức Wikipedia** đầu tiên cho tiếng Việt <sup>70</sup>. Hỗ trợ nghiên cứu các hệ thống kiểm chứng tự động, đặc biệt về **truy hồi bằng chứng** và suy luận phân loại đúng/sai. ViWikiFC được dùng làm bài toán thách thức: hệ thống SemViQA (thắng UIT Workshop) đạt 80,82% độ chính xác nghiêm ngặt, còn mô hình baseline InfoXLM+BM25 chỉ ~67% <sup>71</sup> <sup>72</sup> – cho thấy yêu cầu mô hình phải hiểu ngữ nghĩa sâu và suy luận phức tạp <sup>73</sup>. Dataset này dự kiến công bố rộng rãi, kỳ vọng thúc đẩy nghiên cứu fact-checking dùng tri thức mở (Wikipedia) cho các ngôn ngữ ít tài nguyên.

Đưa ra **benchmark kiểm chứng tin tức đa lĩnh vực** đầu tiên cho tiếng Việt (AAAI 2025) <sup>74</sup> <sup>75</sup>. Mục tiêu đáp ứng nhu cầu cấp bách về công cụ kiểm chứng tự động trong bối cảnh tin giả lan nhanh <sup>76</sup>, nâng độ chính xác của hệ thống kiểm tin tiếng Việt <sup>77</sup>. Nhóm tác giả đã thử nghiệm nhiều mô hình SOTA (kể cả model ngôn ngữ lớn) trên ViFactCheck; đáng chú ý mô hình Gemma đạt F1 macro 89,9%, thiết lập SOTA mới <sup>78</sup> <sup>79</sup>. Dataset kỳ vọng được dùng rộng rãi để huấn luyện & đánh giá hệ thống phát hiện tin sai, và có thể

Làm **tập huấn luyện & đánh giá chính** cho cuộc thi **UIT Data Science Challenge 2023 (task Fact-Checking)** <sup>83</sup> <sup>84</sup>, nơi nhiều đội phát triển mô hình kiểm chứng tự động. Kết quả cuộc thi cho thấy dataset này đã **đẩy mạnh hiệu quả mô hình**: hệ thống thắng cuộc (SemViQA) đạt 78,97% accuracy nghiêm ngặt trên test – vượt xa baseline ban đầu <sup>85</sup>. Sau cuộc thi, ISE-DSC01 tiếp tục được sử dụng trong nghiên cứu học thuật về fact-checking tiếng Việt (huấn luyện BERT/ RoBERTa, thử nghiệm

Tiêu chí	ViAdverNLI (R1–R3)	ViNLI	ViWikiFC	ViFactCheck	ISE-DSC01
	nâng cao độ tin cậy của hệ thống suy luận tự động.			làm nền tảng cho các cuộc thi, ứng dụng thực tiễn trong tương lai <sup>80</sup> <sup>81</sup> . (Dữ liệu và mã nguồn đã công khai trên GitHub <sup>82</sup> ).	phương pháp truy hồi chứng cứ mới) <sup>86</sup> . Nhìn chung, ISE-DSC01 đánh dấu bước tiến lớn, tạo nền tảng dữ liệu lớn thực tế để phát triển & benchmark hệ thống kiểm chứng tự động tiếng Việt <sup>87</sup> .

(Nguồn: tổng hợp thông tin từ các công bố ViNLI <sup>3</sup> <sup>19</sup> <sup>20</sup>, ViWikiFC <sup>13</sup> <sup>62</sup>, ViFactCheck <sup>15</sup> <sup>56</sup>, ISE-DSC01 <sup>7</sup> <sup>40</sup> và ViAdverNLI <sup>2</sup> <sup>18</sup>.)

## Nhận xét so sánh và điểm mạnh của ViAdverNLI

Qua bảng trên, có thể thấy **ViAdverNLI** sở hữu nhiều điểm khác biệt nổi bật so với các bộ dữ liệu trước đây:

- **Độ khó và tính thách thức cao:** ViAdverNLI được thiết kế có chủ đích để *làm khó mô hình*. Điều này thể hiện ở việc mô hình mạnh nhất chỉ đạt ~48% chính xác trên dữ liệu này <sup>53</sup> – thấp hơn đáng kể so với độ chính xác ~79–90% trên ViFactCheck hay ISE-DSC01 (và ~80%+ trên ViWikiFC, ViNLI). ViAdverNLI liên tục *phơi bày điểm yếu* của mô hình qua các vòng, buộc mô hình phải cải thiện khả năng hiểu ngữ cảnh và tránh suy luận dựa trên mẹo. Đây là ưu điểm quan trọng: dataset càng khó sẽ càng thúc đẩy nghiên cứu mô hình NLI **robust** hơn trước những trường hợp đánh lừa.
- **Chiến lược *adversarial* nhiều vòng độc nhất:** Khác với các bộ dữ liệu còn lại xây dựng một lần, ViAdverNLI áp dụng quy trình *human-in-the-loop* qua 3 vòng (R1→R2→R3). Ở mỗi vòng, annotator tìm cách tạo câu hypothesis mà mô hình hiện tại dự đoán sai; những mẫu “đánh bại” được mô hình xác nhận và thêm vào dữ liệu <sup>51</sup>. Chiến lược này giúp ViAdverNLI thu thập được nhiều hiện tượng suy luận **hiếm hóc** (từ việc chơi chữ, tráo đổi từ đồng nghĩa đến tạo bẫy ngữ nghĩa) – những thứ không xuất hiện nhiều trong các dataset chuẩn. Nhờ đó, ViAdverNLI có độ đa dạng cao về cách “bẫy” mô hình, làm cho nó trở thành thước đo tốt để đánh giá khả năng **tổng quát hóa** và chống lại mẫu adversarial của mô hình NLI.

- **Tính đa dạng ngôn ngữ và hiện tượng suy luận:** ViAdverNLI khuyến khích annotator **diễn đạt lại** premise bằng nhiều cách mới (giữ nghĩa hoặc đảo nghĩa) thay vì sao chép từ ngữ. Điều này tạo nên tập dữ liệu với tỷ lệ trùng từ rất thấp và nhiều từ vựng mới trong hypothesis <sup>45</sup>. Các mẫu trong ViAdverNLI bao phủ nhiều hiện tượng: từ phủ định gián tiếp, ẩn dụ, đến thay đổi chi tiết nhỏ (như “30 phút” vs “nửa giờ”) khiến mô hình dễ nhầm lẫn <sup>88</sup> <sup>89</sup>. So với ViNLI hay ViFactCheck – nơi câu hypothesis/claim thường dùng từ ngữ khá sát với premise/bằng chứng – ViAdverNLI *đa dạng và khó đoán* hơn, buộc mô hình học cách **hiểu ngữ nghĩa thực sự** thay vì dựa vào từ khóa đơn thuần.
- **Hiệu quả trong cải thiện mô hình:** Một lợi thế đáng chú ý là khi huấn luyện mô hình trên ViAdverNLI, hiệu suất trên các bộ dữ liệu khác được cải thiện rõ rệt <sup>53</sup>. Điều này gợi ý rằng ViAdverNLI không chỉ đóng vai trò kiểm tra độ bền của mô hình, mà còn như một bộ dữ liệu huấn luyện giúp mô hình **tổng quát tốt hơn** (có lẽ do học được cách xử lý các bẫy ngôn ngữ). Ngược lại, các dataset khác tuy hữu ích để đánh giá thông thường, nhưng có thể chưa đủ phong phú để nâng cao độ bền cho mô hình trước các ví dụ “hiểm”.
- **Bổ sung cho khoảng trống tài nguyên:** Trước ViAdverNLI, các benchmark NLI tiếng Việt (như ViNLI, VnNewsNLI) đều tập trung vào dữ liệu “thẳng” và chất lượng, chưa chú trọng khía cạnh adversarial <sup>90</sup>. ViAdverNLI ra đời đã **bổ khuyết khoảng trống** này, đưa tiếng Việt vào xu hướng nghiên cứu NLI tiên tiến giống như tiếng Anh (ANLI) hay tiếng Trung <sup>91</sup>. Trong bối cảnh phát hiện sai lệch thông tin ngày càng quan trọng, dataset này đóng góp cách đánh giá mới, khuyến khích phát triển mô hình NLI và fact-checking có khả năng chống chịu tốt hơn với thông tin gây nhiễu.

Tóm lại, **ViAdverNLI** nổi trội nhờ cách xây dựng sáng tạo và độ phức tạp cao, giúp kiểm tra và huấn luyện mô hình ở một mức độ **khó** hơn so với các bộ dữ liệu hiện có. Trong khi ViNLI cung cấp nền tảng cơ bản về suy luận, ViWikiFC mở rộng sang tri thức Wikipedia, ISE-DSC01 và ViFactCheck đem lại khối lượng dữ liệu lớn và đa miền cho fact-checking, thì ViAdverNLI tập trung vào khía cạnh *thử thách mô hình*, khiến nó trở thành **benchmark đặc biệt giá trị** để thúc đẩy các nghiên cứu NLI/fact-checking tiếng Việt lên tầm cao mới. Các điểm mạnh này của ViAdverNLI bổ sung cho những benchmark sẵn có, cùng nhau tạo nên một bộ dữ liệu đa dạng giúp phát triển các hệ thống kiểm chứng thông tin và suy luận ngôn ngữ tự nhiên ngày càng toàn diện và đáng tin cậy hơn.

**Nguồn tài liệu tham khảo:** ViNLI <sup>3</sup> <sup>20</sup>; ViWikiFC <sup>13</sup> <sup>71</sup>; ViFactCheck <sup>15</sup> <sup>56</sup>; ISE-DSC01 <sup>7</sup> <sup>40</sup>; ViAdverNLI <sup>18</sup> <sup>53</sup>.

<sup>1</sup> <sup>2</sup> <sup>8</sup> <sup>9</sup> <sup>17</sup> <sup>18</sup> <sup>41</sup> <sup>42</sup> <sup>45</sup> <sup>46</sup> <sup>47</sup> <sup>51</sup> <sup>52</sup> <sup>53</sup> <sup>57</sup> <sup>88</sup> <sup>89</sup> <sup>90</sup> <sup>91</sup> ViANLI: Adversarial Natural Language Inference for Vietnamese  
<https://arxiv.org/html/2406.17716v2>

<sup>3</sup> <sup>4</sup> <sup>5</sup> <sup>6</sup> <sup>7</sup> <sup>13</sup> <sup>14</sup> <sup>15</sup> <sup>16</sup> <sup>19</sup> <sup>20</sup> <sup>21</sup> <sup>22</sup> <sup>23</sup> <sup>24</sup> <sup>25</sup> <sup>26</sup> <sup>27</sup> <sup>28</sup> <sup>29</sup> <sup>30</sup> <sup>31</sup> <sup>32</sup> <sup>33</sup> <sup>34</sup> <sup>35</sup> <sup>36</sup> <sup>37</sup> <sup>38</sup> <sup>39</sup>  
<sup>40</sup> <sup>54</sup> <sup>55</sup> <sup>56</sup> <sup>62</sup> <sup>63</sup> <sup>64</sup> <sup>65</sup> <sup>66</sup> <sup>67</sup> <sup>68</sup> <sup>69</sup> <sup>70</sup> <sup>71</sup> <sup>72</sup> <sup>73</sup> <sup>74</sup> <sup>75</sup> <sup>76</sup> <sup>77</sup> <sup>78</sup> <sup>79</sup> <sup>80</sup> <sup>81</sup> <sup>82</sup> <sup>83</sup> <sup>84</sup> <sup>85</sup> <sup>86</sup> <sup>87</sup>  
 VNese Fact-Checking Data Collection.pdf  
<file://file-UxTqzMGCQumdBMqzFR2UFu>

<sup>10</sup> <sup>11</sup> <sup>12</sup> <sup>43</sup> <sup>44</sup> <sup>48</sup> <sup>49</sup> <sup>50</sup> <sup>58</sup> <sup>59</sup> <sup>60</sup> <sup>61</sup> ViNLI: A Vietnamese Corpus for Studies on Open-Domain Natural Language Inference  
<https://aclanthology.org/2022.coling-1.339.pdf>