

AI-Assisted Competency Assessment from Egocentric Video in Simulation-Based Nursing Education

Anonymous CVPR submission

Paper ID 18

Abstract

001 *Assessing learner competency in clinical simulation re-*
002 *quires expert observation that is time-intensive, difficult to*
003 *scale, and subject to inter-rater variability. While mul-*
004 *timodal learning analytics has emerged as a promising*
005 *direction, the contributions of individual modalities re-*
006 *main underexplored. We investigate what vision perspec-*
007 *tive alone can reveal by proposing a three-stage frame-*
008 *work that (1) extracts action timelines from egocentric*
009 *nursing simulation video using frozen visual encoders and*
010 *few-shot learning, (2) derives sequence-level features and*
011 *per-session recognition metrics, and (3) relates these to*
012 *instructor-rated competency. Across 22 densely annotated*
013 *sessions (3.8 hours, 493 actions), a frozen DINOv2 back-*
014 *bone with HMM Viterbi decoding achieves 57.4% MOF in*
015 *leave-one-out 1-shot recognition. Surprisingly, we observe*
016 *a negative trend between recognition accuracy and compe-*
017 *tency ($\rho = -0.524$, $p = 0.012$ for $mIoU$), robust to six con-*
018 *found controls: more competent students produce diverse,*
019 *harder-to-classify workflows, while simple sequence fea-*
020 *tures show no such relationship. Per-item analysis identifies*
021 *patient safety protocols and team communication as the ex-*
022 *pected behaviors most reflected in this pattern, and process*
023 *model comparisons reveal that high-competency students*
024 *exhibit more protocol-consistent action transitions. These*
025 *findings suggest that recognition accuracy may complement*
026 *predicted action timelines as a pedagogically informative*
027 *signal in automated competency assessment.*

028 1. Introduction

029 Across education and workforce training, a central goal is
030 to determine whether learners have developed the knowl-
031 edge, skills, and judgment needed to perform effectively
032 in practice, a quality broadly termed *competency* [32, 46].
033 In domains defined by skilled physical performance, com-
034 petency assessment requires expert observation of context-
035 dependent behaviors that unfold over time [10]. The conse-

quences of undetected gaps are especially acute in clinical 036
education, where medication administration errors remain 037
among the most common preventable adverse events, of- 038
ten rooted in procedural lapses missed during training [17]. 039
Simulation-based learning addresses this by letting students 040
practice clinical skills without risk to real patients [17, 22], 041
but competency encompasses not just executing procedures 042
correctly but doing so in an appropriate sequence, with 043
complete safety checks and fluid task transitions [10, 26]. 044
Instructors assess each session using standardized instru- 045
ments such as the Creighton Competency Evaluation Instru- 046
ment (C-CEI) [44] that map observable behaviors to compe- 047
tency constructs (e.g., clinical judgment, patient safety) [1, 048
26]. This model faces two structural constraints: expert ob- 049
servation cannot scale with growing cohorts [3], and inter- 050
rater reliability remains only moderate to substantial even 051
among trained faculty [20, 23]. Multimodal Learning Ana- 052
lytics (MMLA) [6, 9] has emerged as a promising response, 053
integrating video, audio, physiological signals, and inter- 054
action logs to study learning processes at scale [2, 37]. 055
However, in the context of nursing simulation competency, 056
the contributions of individual modalities remain underex- 057
plored. As a first step, we focus on vision [43] to investigate 058
what this modality can reveal. 059

Video captures what learners do, how they move, and 060
what objects they interact with, all without requiring in- 061
strumented environments. Recent advances in first-person 062
video understanding have made egocentric recordings espe- 063
cially compelling [15]. Head-mounted cameras provide 064
an unobstructed, hands-proximal view of what a student at- 065
tends to and acts upon, and when paired with gaze sens- 066
ing [7, 21], can reveal attentional patterns linked to errors in 067
skilled activities [27]. As shown in Fig. 1, this perspective 068
captures the full behaviors of a clinical encounter, from how 069
students hold a medication bottle to which device they use 070
for dosage calculation, preserving precisely the signatures 071
that distinguish novice from expert workflows. 072

Medication administration is an ideal proving ground be- 073
cause competence is inherently sequential: correct actions 074
in the wrong order, or with safety steps omitted, consti- 075

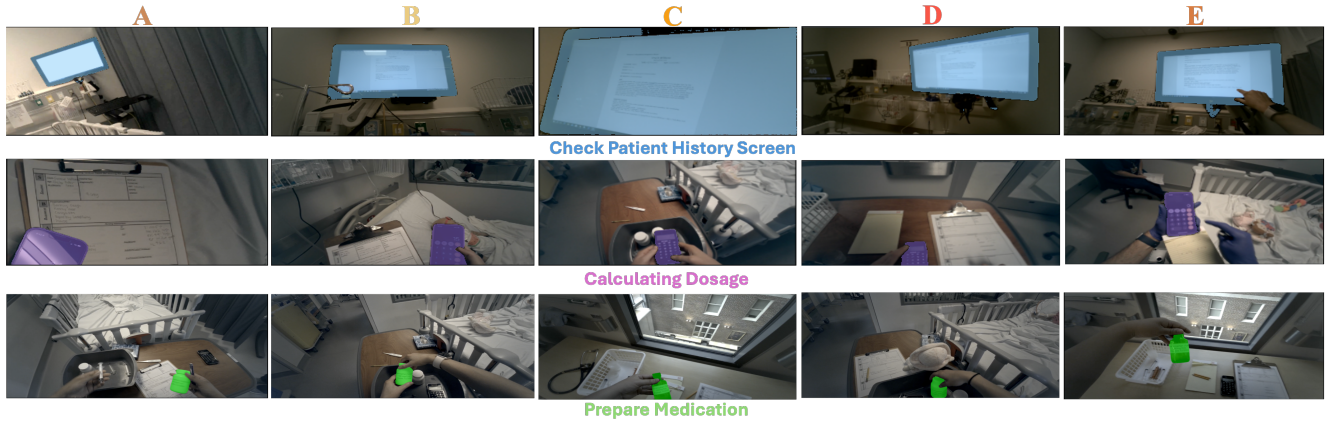


Figure 1. Example images of [checking the patient screen](#), [calculating dosage](#), and [preparing medication](#) from simulation videos of five nursing students. Students A, B, and E use phones for dosage calculation, whereas students C and D use handheld calculators. During medication preparation, students C and E use a dark brown medicine bottle, while students A, B, and D each use a different bottle and hold it differently. More details of the simulation procedure are in App. A.

076 tute a clinical error [22]. The annotation codebook used in
 077 this study was developed around the medication administra-
 078 tion workflow, drawing on established nursing competency
 079 measurement tools [36], and captures fine-grained actions
 080 such as dosage calculation. The C-CEI rubric maps ex-
 081 pected behaviors to broader competency constructs across
 082 the full clinical encounter; vision-based analysis can speak
 083 directly to the former, while the latter requires instructor
 084 judgment. Applying egocentric video understanding to this
 085 setting, however, introduces domain-specific challenges:
 086 clinical action vocabularies are absent from standard bench-
 087 marks [5, 8], cohorts are small due to privacy constraints,
 088 and models pre-trained on real bodies face a domain gap
 089 with simulation mannequins. While surgical education has
 090 demonstrated that AI can recognize operative phases and
 091 classify skill levels from laparoscopic recordings [24], that
 092 work assumes data-rich, fixed-camera environments. Nurs-
 093 ing simulation differs in that the recording is egocentric, co-
 094 horts are governed by IRB constraints, and competency is
 095 holistic rather than tied to a single technical procedure. Yet
 096 these difficulties may not be purely noise. Research in sur-
 097 gical skill assessment has found that automated classifiers
 098 perform worse on higher-skilled practitioners [41], and that
 099 temporal patterns of action execution capture skill level bet-
 100 ter than outcome measures alone [47]. This raises the possi-
 101 bility that recognition accuracy itself carries a pedagogical
 102 signal, with lower accuracy reflecting the diverse workflows
 103 of more competent students.

104 This gap motivates the present work to investigate
 105 whether egocentric video, analyzed through frozen visual
 106 encoders and few-shot learning, can serve as the basis for
 107 automated competency assessment in nursing simulation.
 108 Specifically, we ask whether action timelines can be re-

liably extracted from first-person clinical video under ex-
 treme low-data constraints; whether the resulting sequence
 features and recognition difficulty relate to instructor-rated
 competency; and whether temporal action patterns differ
 between high- and low-performing students. To answer
 these questions, we propose a three-stage framework that
 extracts action timelines, analyzes their sequential structure,
 and relates sequence features and recognition accuracy to
 competency scores across 22 densely annotated sessions.

Research Questions.

- RQ1.** To what extent can few-shot action recognition iden-
 tify clinical actions in egocentric nursing simulation
 video? (§5.1)
- RQ2.** To what extent do automatically extracted ac-
 tion sequences and recognition difficulty relate to
 instructor-rated competency, and which expected be-
 haviors on the C-CEI are most reflected in vision-
 based action analysis? (§5.2)
- RQ3.** What temporal action patterns distinguish high- from
 low-performing students? (§5.3)

Contributions. This work advances competency assess-
 ment from egocentric video through three contributions:

- Few-shot clinical action recognition.** We show that frozen DINOv2 features with HMM Viterbi decoding achieve 57.4% MOF in leave-one-out 1-shot action recognition of egocentric nursing simulation video, establishing feasibility under extremely low-data conditions without any fine-tuning.
- Classification difficulty to competency.** We observe a negative trend between recognition accuracy and instructor-rated competency ($\rho = -0.524$, $p = 0.012$ for mIoU), robust to six confound controls. Per-item analysis identifies expected behaviors related to patient

142 safety protocols and team communication as the C-CEI
143 items most reflected in this pattern.

144 3. **Temporal workflow analysis.** Process model compar-
145 isons from ground-truth action sequences show
146 that high-competency students exhibit more diverse,
147 protocol-consistent action transitions, delineating the
148 boundaries of unimodal video-based assessment.

149 2. Related Work

150 **Computer vision for clinical skill assessment.** Deep
151 learning has been applied extensively to surgical work-
152 flow recognition and skill evaluation from operative video,
153 including tool detection and phase recognition [45], di-
154 rect skill classification from video [13], fine-grained ac-
155 tion triplet recognition [29], and scalable objective assess-
156 ment of technical skill [14, 16]. These studies highlight
157 the promise of video-based clinical assessment, but most
158 are developed in data-rich surgical settings with fixed cam-
159 eras and relatively controlled workflows. In contrast, our
160 setting is egocentric and small-scale, and the goal is to as-
161 sess holistic nursing competency rather than isolated tech-
162 nical skill. **Computer vision for education and learn-
163 ing analytics.** MMLA integrates video, audio, physio-
164 logical signals, and interaction logs to study learning pro-
165 cesses [30]. Within this paradigm, vision has been used
166 to detect learning-relevant affective states [4], align neural
167 attention with human gaze [42], analyze embodied class-
168 room learning [11, 12], and model student interaction se-
169 quences [40]; a recent review surveys multimodal meth-
170 ods across adult training environments, including nursing
171 simulation [38]. We focus on a single modality (egocentric
172 video) to establish what vision alone can reveal about clin-
173 ical competency. **Few-shot and temporal action recogni-
174 tion.** Prototype networks enable classification with mini-
175 mal labeled examples. Temporal action segmentation has
176 advanced rapidly on standard benchmarks [8], and large-
177 scale egocentric datasets [15] together with self-supervised
178 encoders such as DINOv2 [31] provide strong frozen rep-
179 resentations. We combine prototype matching with HMM
180 Viterbi decoding [33] for temporal segmentation under ex-
181 treme low-data clinical conditions.

182 3. Problem Formulation

183 Let $\mathcal{V} = \{V_1, \dots, V_N\}$ denote a set of N egocentric video
184 sessions, where each session $V_i = (f_i^{(1)}, f_i^{(2)}, \dots, f_i^{(T_i)})$
185 consists of T_i ordered frames, where $i \in \{1, \dots, N\}$ in-
186 dexes the session and $t \in \{1, \dots, T_i\}$ indexes the frame.
187 Each video is associated with an instructor-assigned com-
188 petency score vector $\mathbf{c}_i \in \mathbb{R}^{23}$ across 23 expected behav-
189 iors on the C-CEI rubric (App. B), of which 11 correspond
190 to video-observable actions; the mean of these 11 items
191 yields the overall competency percentage used for associ-

192 ation analyses. Not all items are rated for every session, so
193 some entries of \mathbf{c}_i are missing.

194 *Stage 1: Action Recognition.* Given a clinically grounded
195 action taxonomy $\mathcal{A} = \{a_1, \dots, a_K, a_\emptyset\}$ comprising $K=16$
196 clinical action classes and one background class a_\emptyset (17 la-
197 bels total), the goal is to assign each frame a label $y_i^{(t)} \in \mathcal{A}$,
198 producing a frame-level prediction $\hat{\mathbf{y}}_i = (\hat{y}_i^{(1)}, \dots, \hat{y}_i^{(t)})$.
199 A frozen encoder ϕ extracts per-frame features $\mathbf{z}_i^{(t)} =$
200 $\phi(f_i^{(t)})$, which are matched against class prototypes com-
201 puted from a support set \mathcal{S} of labeled exemplars sampled
202 from held-out sessions:

$$203 \hat{\mathbf{y}}_i = \text{Decode}(\text{sim}(\phi(V_i), \mathcal{P}(\mathcal{S}))), \quad (1)$$

204 where $\mathcal{P}(\mathcal{S})$ computes class prototypes from the support
205 set [39] (Sec. 4.1) and Decode applies HMM Viterbi de-
206 coding [33] to enforce temporally coherent label sequences.
207 We evaluate $\hat{\mathbf{y}}_i$ against ground-truth annotations \mathbf{y}_i^* using
208 frame-level accuracy (MOF), mean intersection-over-union
209 (mIoU), and segmental F1 (RQ1).

210 *Stage 2: Sequence Analysis.* From the predicted
211 frame-level labels $\hat{\mathbf{y}}_i$, we collapse contiguous same-label
212 frames into an ordered action sequence $\mathbf{s}_i =$
213 $((c_i^{(1)}, d_i^{(1)}), \dots, (c_i^{(L_i)}, d_i^{(L_i)}))$ of L_i segments ($l \in$
214 $\{1, \dots, L_i\}$), where $c_i^{(l)} \in \mathcal{A} \setminus \{a_\emptyset\}$ is the action label and
215 $d_i^{(l)}$ is the segment duration in frames. From this sequence
216 we derive two families of features used in subsequent anal-
217 yses: (1) action transition frequencies, which capture the
218 pairwise flow between clinical actions, and (2) per-video
219 recognition metrics (MOF, mIoU, F1), which summarize
220 how well the classifier fits each session.

221 *Stage 3: Competency Analysis.* Given the small sample
222 size ($N = 22$) and the pedagogical requirement for trans-
223 parent feedback, we map sequence features to competency
224 scores using Spearman rank association. To disentangle
225 action detection errors from the intrinsic limits of vision-
226 based assessment, we evaluate under both oracle (features
227 from ground-truth \mathbf{y}_i^*) and predicted (features from $\hat{\mathbf{y}}_i$) con-
228 ditions. Per-item analysis examines which expected be-
229 haviors are captured by action sequences alone, and com-
230 parison of ground-truth action transition graphs across per-
231 formance groups identifies discriminative temporal patterns
232 (RQ2, RQ3).

233 4. Method

234 We collect egocentric video from 22 nursing students, each
235 performing a single standardized pediatric simulation on
236 high-fidelity mannequins. Each session captures one stu-
237 dent’s complete first-person view of the clinical encounter,
238 recorded via egocentric glasses at 25 FPS. Sessions range
239 from 4 to 24 minutes (mean 10.5 min, total 3.8 hours).

240 Each session is annotated across three temporal layers by

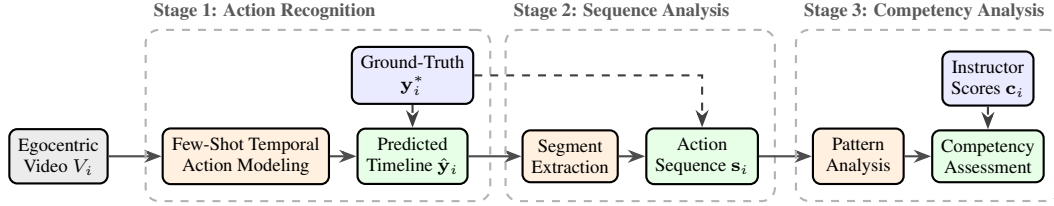


Figure 2. Overview of the proposed three-stage framework. Gray boxes denote inputs, orange boxes denote processing modules, green boxes denote outputs, and blue boxes denote supervision or reference signals. Solid arrows indicate the forward inference flow, while dashed arrows indicate supervision or oracle guidance. The three stages perform action timeline prediction, action sequence construction, and competency assessment, respectively.

241 a trained coder using the NOVA annotation system: (1) Behaviors (3 classes: Introduction, Assessment, Administration), (2) Actions ($K = 16$ fine-grained clinical classes plus one background class a_{\emptyset} for unannotated frames, yielding 17 labels total; see App. C), and (3) Communication (Patient, Family, Provider). The Action layer, containing 493 annotated clinical segments, serves as the primary target for few-shot recognition. Each session is independently rated by an expert instructor across 23 expected behaviors using the C-CEI rubric (App. B), of which the 11 video-observable items yield the overall competency percentage used throughout this study. Inter-rater reliability was assessed on 3 stratified videos (low/median/high competency) independently annotated by a second rater, yielding substantial agreement (mean Cohen’s $\kappa = 0.708$; App. F).

256 4.1. Stage 1: Action Recognition

257 **Feature extraction.** We extract frame-level features from each video using a frozen backbone encoder. For each frame $f_i^{(t)}$, we obtain a feature vector $\mathbf{z}_i^{(t)} = \phi(f_i^{(t)}) \in \mathbb{R}^D$, followed by L2 normalization: $\mathbf{z}_i^{(t)} \leftarrow \mathbf{z}_i^{(t)} / \|\mathbf{z}_i^{(t)}\|$. We evaluate three backbones: (1) ResNet-50 [18] (ImageNet-supervised, $D = 2048$), (2) DINOv2 ViT-B/14 [31] (ImageNet self-supervised, $D = 768$), and (3) CLIP ViT-B/16 [34] (vision-language contrastive, $D = 512$). All backbones are frozen with no fine-tuning.

266 **Prototype computation.** In the cross-sample (leave-one-out) setting, we construct class prototypes from the $N-1$ support sessions following the prototypical network paradigm [39]. For each support session V_j and each class k present in that session, we randomly sample n labeled frames and compute a per-session centroid $\boldsymbol{\mu}_{k,j}$, which is then L2-normalized. Because not every session contains every action class, the global prototype for class k is obtained by averaging the normalized centroids only over sessions that contain that class: $\mathbf{p}_k = \frac{1}{|\mathcal{J}_k|} \sum_{j \in \mathcal{J}_k} \frac{\boldsymbol{\mu}_{k,j}}{\|\boldsymbol{\mu}_{k,j}\|}$, where \mathcal{J}_k is the set of sessions containing class k . The per-session normalization ensures each session contributes a unit-direction vector, preventing sessions whose sampled frames are more self-consistent from dominating the pro-

280 totype direction. The aggregated prototype is then L2-normalized again to ensure it lies on the unit sphere, since the mean of unit vectors is not itself unit-length in general. As an alternative, we also evaluate a clustered strategy in which all support frames for each class are pooled and partitioned into $k=3$ sub-centroids via k -means; each query frame is then assigned to the class of its nearest sub-centroid.

288 **Classification.** Each query frame is scored against all prototypes (including the background class a_{\emptyset}) via cosine similarity. Rather than committing to a hard per-frame label at this stage, the continuous similarity scores are passed directly to the temporal smoothing step below, which jointly optimizes over the entire sequence.

294 **Temporal smoothing.** We apply *HMM Viterbi decoding* [33] to enforce temporally coherent label sequences. The transition matrix \mathbf{A} is learned from support session label sequences with Laplace smoothing, prior probabilities $\boldsymbol{\pi}$ are estimated from action start frequencies, and emission log-probabilities are computed as temperature-scaled ($\tau = 5$, selected via grid search on held-out folds; values in the range 5–10 are standard for cosine-prototype matching in few-shot settings [39]) log-softmax of cosine similarities:

$$293 \log P(\mathbf{z}_i^{(t)} | a_k) = \tau \cdot \cos(\mathbf{z}_i^{(t)}, \mathbf{p}_k) - \log \sum_{k'} \exp(\tau \cdot \cos(\mathbf{z}_i^{(t)}, \mathbf{p}_{k'})). \quad (2)$$

304 The Viterbi algorithm [33] then recovers the optimal label sequence $\hat{\mathbf{y}}_i = (\hat{y}_i^{(1)}, \dots, \hat{y}_i^{(T_i)})$ by selecting, at each time step, the action label that jointly maximizes the cumulative sum of emission log-probabilities and transition log-probabilities over the entire sequence, thereby enforcing clinically plausible transitions rather than treating each frame independently. This smoothed timeline is the final output of Stage 1 (Sec. 4.1) and the prediction used in all subsequent analyses (Sec. 4.2–4.3).

313	4.2. Stage 2: Sequence Analysis	
314	From the frame-level predictions \hat{y}_i , we collapse contiguous	362
315	same-label frames into an ordered sequence of action	363
316	segments via run-length encoding, discarding segments be-	364
317	low a minimum duration threshold and removing back-	365
318	ground segments. From the resulting clinical action se-	366
319	quence we compute two families of features. First, we tab-	367
320	ulate pairwise action transition frequencies, which record	368
321	how often each action class is followed by every other class	369
322	within a session; these frequencies form the basis of the pro-	370
323	cess model comparison in Sec. 5.3. Second, we retain the	371
324	per-video frame-level recognition metrics (MOF, mIoU, F1)	372
325	computed during Stage 1, which serve as summary mea-	373
326	sures of how well the classifier fits each session and are used	
327	in the competency analysis (Sec. 5.2).	
328	4.3. Stage 3: Competency Analysis	
329	Given the small sample size and the pedagogical need for	
330	transparency, we employ Spearman’s rank correlation to	
331	test the relationship between sequence-level features and	
332	video-observable competency scores (11 items). To assess	
333	robustness, we compute partial associations controlling for	
334	potential confounds (annotation coverage, video duration,	
335	segment count). Per-item analysis examines which of the	
336	23 expected behaviors are most reflected in recognition ac-	
337	curacy. Process model analysis (Heuristics Miner) visual-	
338	izes action transition graphs for high- and low-competency	
339	groups to identify differences in clinical workflows.	
340	5. Results	
341	5.1. RQ1: Few-Shot Clinical Action Recognition	
342	We evaluate few-shot action recognition under two set-	
343	tings: <i>within-sample</i> , where support and query frames are	
344	drawn from the same video, and <i>cross-sample</i> (leave-one-	
345	out), where the model must generalize to entirely unseen	
346	sessions. We report three standard temporal action segmen-	
347	tation metrics: mean-over-frames accuracy (MOF), mean	
348	intersection-over-union (mIoU), and segmental F1 score.	
349	5.1.1. Within-Sample Evaluation	
350	In the within-sample setting, for each of the 22 videos, n	
351	frames per action class are randomly sampled as support	
352	prototypes (where n denotes the shot count); the remain-	
353	ing frames serve as the query set. Tab. 1 reports within-	
354	sample performance across five shot counts. Recognition	
355	quality improves substantially with more support examples,	
356	with the largest gain between 1 and 3 shots. Performance	
357	plateaus around 10–15 shots, indicating that even a mod-	
358	est number of labeled exemplars enables reliable within-	
359	sample segmentation and that the bottleneck lies in cross-	
360	session generalization rather than representation capacity.	
	5.1.2. Cross-Sample Evaluation	361
	The more challenging and practically relevant setting is	362
	cross-sample evaluation, where the model must generalize	363
	to entirely unseen sessions with unseen participants. We	364
	adopt a leave-one-out protocol across all 22 sessions: in	365
	each of 22 folds, one session is held out as the query video,	366
	and class prototypes are constructed from the remaining 21	367
	support sessions using prototype computation (Sec. 4.1).	368
	We vary the shot count $n \in \{1, 3, 5, 10, 15\}$ frames sam-	369
	pled per class per session to examine how prototype quality	370
	scales with the support budget. The query video is classified	371
	frame-by-frame via cosine similarity followed by HMM	372
	Viterbi decoding.	373
	Tab. 2 reports cross-sample performance for DINOv2,	374
	ResNet-50 and CLIP across five shot counts (1, 3, 5,	375
	10, 15), comparing mean and clustered prototype strate-	376
	gies. DINOv2 with mean prototypes consistently outper-	377
	forms all other configurations, achieving its best perfor-	378
	mance at 10 shots (65.6% MOF, 45.1% mIoU, 41.9% F1).	379
	Mean prototypes substantially outperform clustered proto-	380
	types across both backbones, indicating that splitting each	381
	class into multiple sub-centroids introduces false matches	382
	under the few-shot budget. DINOv2 consistently outper-	383
	forms ResNet-50 and CLIP across all shot counts and met-	384
	rics, suggesting that self-supervised vision transformer fea-	385
	tures offer better discrimination of fine-grained hand-object	386
	interactions in clinical settings.	387
	Comparing Tab. 2 with Tab. 1, the cross-sample set-	388
	ting shows substantially lower performance than within-	389
	sample at the same nominal shot count n . Importantly, these	390
	are not directly comparable in terms of total support data:	391
	within-sample uses n frames per class from a single video,	392
	whereas cross-sample pools n frames per class from each	393
	of 21 sessions, yielding 20 times more support frames per	394
	class overall. Despite this $21 \times$ larger support budget, cross-	395
	sample performance lags behind, underscoring that the gap	396
	is driven by participant-level domain shift (differences in	397
	student appearance, camera angle, pace, and workflow or-	398
	dering across individuals) rather than insufficient support	399
	data. Yet the following sections show that this recognition	400
	variability is itself pedagogically informative.	401
	5.2. RQ2: Action Sequences, Competency, and Per-	402
	Item Analysis	403
	To investigate whether extracted action recognition metrics	404
	carry information related to instructor-rated competency,	405
	we compute Spearman rank associations between per-video	406
	recognition metrics obtained from cross-sample (leave-one-	407
	out) evaluation using the best model (DINOv2 + HMM,	408
	10-shot) and the overall competency score (mean of the 11	409
	video-observable rubric items; see Sec. 3).	410

Table 1. Within-sample few-shot action recognition. For each video, n frames per action class are sampled as prototypes, and the remaining frames serve as the query set. Results are reported as mean \pm std over 22 videos. For all backbones, **bold** indicates the best-performing configuration per metric and underlined indicates the second-best. Higher is better for all metrics.

Shots (n)	DINOv2			ResNet50			CLIP (ViT-B/16)		
	MOF	mIoU	F1	MOF	mIoU	F1	MOF	mIoU	F1
1	0.555 \pm 0.125	0.419 \pm 0.108	0.522 \pm 0.115	0.559 \pm 0.150	0.416 \pm 0.130	0.514 \pm 0.140	0.556 \pm 0.135	0.393 \pm 0.113	0.495 \pm 0.124
3	0.783 \pm 0.083	0.632 \pm 0.088	0.736 \pm 0.077	0.786 \pm 0.075	0.637 \pm 0.078	0.737 \pm 0.072	0.797 \pm 0.081	0.642 \pm 0.086	0.744 \pm 0.076
5	0.847 \pm 0.072	0.715 \pm 0.099	0.805 \pm 0.076	0.836 \pm 0.079	0.699 \pm 0.106	0.791 \pm 0.084	0.828 \pm 0.082	0.685 \pm 0.104	0.779 \pm 0.086
10	0.905 \pm 0.054	0.797 \pm 0.081	0.852 \pm 0.067	0.896 \pm 0.056	0.795 \pm 0.074	0.852 \pm 0.065	0.908 \pm 0.041	0.798 \pm 0.069	0.850 \pm 0.067
15	0.930 \pm 0.042	0.846 \pm 0.077	0.869 \pm 0.078	0.923 \pm 0.044	0.833 \pm 0.090	0.856 \pm 0.090	0.931 \pm 0.030	0.844 \pm 0.066	0.858 \pm 0.064

Table 2. Cross-sample few-shot action recognition (leave-one-out, 22 folds). For each held-out video, n frames per action class are sampled from the 21 support sessions to construct prototypes, and the held-out session is classified via HMM Viterbi decoding. Results are reported as mean \pm std over 22 folds. **bold** indicates the best-performing configuration per metric within each backbone, and underlined indicates the second-best. Higher is better for all metrics.

Shots (n)	Proto.	DINOv2			ResNet50			CLIP (ViT-B/16)		
		MOF	mIoU	F1	MOF	mIoU	F1	MOF	mIoU	F1
1	Mean	0.574 \pm 0.121	0.337 \pm 0.088	0.341 \pm 0.092	0.467 \pm 0.139	0.255 \pm 0.083	0.256 \pm 0.073	0.533 \pm 0.134	0.306 \pm 0.098	0.322 \pm 0.095
	Clust.	0.472 \pm 0.111	0.280 \pm 0.088	0.245 \pm 0.061	0.397 \pm 0.126	0.225 \pm 0.084	0.195 \pm 0.055	0.464 \pm 0.160	0.292 \pm 0.105	0.255 \pm 0.073
3	Mean	0.622 \pm 0.128	0.433 \pm 0.122	0.412 \pm 0.106	0.507 \pm 0.113	0.302 \pm 0.078	0.287 \pm 0.058	0.580 \pm 0.137	0.405 \pm 0.132	0.373 \pm 0.111
	Clust.	0.569 \pm 0.104	0.403 \pm 0.103	0.342 \pm 0.084	0.445 \pm 0.106	0.257 \pm 0.064	0.226 \pm 0.063	0.502 \pm 0.112	0.337 \pm 0.108	0.279 \pm 0.081
5	Mean	0.640 \pm 0.119	0.436 \pm 0.109	0.411 \pm 0.096	0.508 \pm 0.139	0.327 \pm 0.111	0.304 \pm 0.087	0.596 \pm 0.115	0.405 \pm 0.099	0.390 \pm 0.099
	Clust.	0.594 \pm 0.134	0.423 \pm 0.115	0.372 \pm 0.107	0.394 \pm 0.110	0.247 \pm 0.086	0.224 \pm 0.061	0.540 \pm 0.136	0.371 \pm 0.098	0.333 \pm 0.083
10	Mean	0.656 \pm 0.152	0.451 \pm 0.128	0.419 \pm 0.100	0.473 \pm 0.138	0.298 \pm 0.085	0.279 \pm 0.076	0.618 \pm 0.124	0.439 \pm 0.098	0.414 \pm 0.104
	Clust.	0.612 \pm 0.131	0.433 \pm 0.117	0.407 \pm 0.112	0.428 \pm 0.123	0.282 \pm 0.101	0.268 \pm 0.089	0.574 \pm 0.116	0.400 \pm 0.084	0.374 \pm 0.076
15	Mean	0.644 \pm 0.148	0.446 \pm 0.128	0.417 \pm 0.106	0.496 \pm 0.163	0.299 \pm 0.116	0.290 \pm 0.109	0.612 \pm 0.134	0.426 \pm 0.097	0.407 \pm 0.088
	Clust.	0.594 \pm 0.128	0.399 \pm 0.106	0.391 \pm 0.099	0.433 \pm 0.146	0.256 \pm 0.090	0.249 \pm 0.078	0.550 \pm 0.149	0.385 \pm 0.096	0.364 \pm 0.087

411 5.2.1. Overall Trends

412 Tab. 3 presents a notable pattern: all three recognition accuracy
413 metrics show a negative trend as video-observable
414 competency increases. The strongest observed relationship
415 is for mIoU ($\rho = -0.524$, $p = 0.012$), which measures
416 per-class balance. MOF ($\rho = -0.439$, $p = 0.041$) and F1
417 ($\rho = -0.433$, $p = 0.044$) show similar patterns. Neither
418 the number of ground-truth action classes nor the number
419 of labeled query frames shows any significant relationship
420 with competency, ruling out annotation-count artifacts as a
421 confounding explanation.

Table 3. Spearman ρ between per-video recognition metrics and overall video-observable competency score (11 items, $N = 22$). All accuracy metrics show a negative trend with competency.

Feature	Spearman ρ	p -value	Pearson r
mIoU (per-class accuracy)	-0.524	0.012	-0.469
MOF (frame accuracy)	-0.439	0.041	-0.436
F1 (macro)	-0.433	0.044	-0.407
Frame error rate (1-MOF)	+0.439	0.041	+0.436
# Action classes in GT	-0.087	0.701	-0.043
# Labeled query frames	-0.071	0.754	+0.016

422 This pattern is consistent with Moravec’s insight [28]:

the classifier performs better on the mechanical, templated
workflows of lower-performing students, whereas the fluid,
adaptive behaviors of higher competence prove harder to
recognize. When sessions are split by median competency
score, LOW-competency students have 9.5% higher MOF
and 8.3% higher mIoU than HIGH-competency students
(Fig. 3). One plausible interpretation, consistent with the
motor learning principle of abundance [35], is that more
competent students perform more diverse workflows with
additional safety checks and fluid task transitions, produc-
ing greater visual diversity that makes classification harder
while earning higher instructor marks. This converges with
surgical skill assessment findings where classifiers achieve
lower accuracy on higher-skilled practitioners [19, 41], sug-
gesting that the negative trend in recognition accuracy *may*
carry a pedagogically informative signal. We note, how-
ever, that with $N = 22$ sessions, this interpretation remains
exploratory. Importantly, simple sequence features (screen
time ratio, transition count, unique action count) extracted
from both oracle and predicted timelines show no signifi-
cant relationship with competency (all $p > 0.10$).

5.2.2. Per-Item Analysis

To identify which facets of clinical competency are most ac-
cessible through vision-based action analysis, we examine

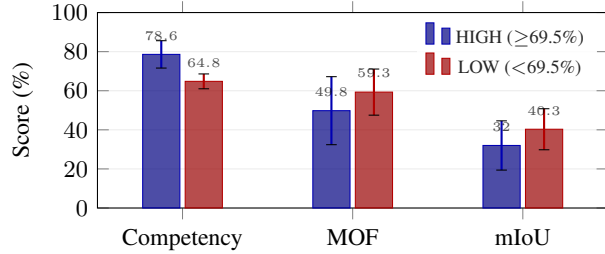


Figure 3. Group-level comparison: sessions split by median video-observable competency score (11 items). Despite higher instructor ratings, HIGH-competency students exhibit *lower* classification accuracy. Error bars show ± 1 std.

447 Spearman associations between per-video MOF and each of
448 the 23 instructor rubric items. Tab. 4 reports the five items
449 with the strongest associations, ordered by magnitude.

Table 4. Top-5 Spearman ρ between per-video MOF and individual rubric items, ordered by magnitude. Shaded rows denote items with visually observable behaviors. N varies across items because instructors may omit ratings when a behavior is not observed or not applicable during a particular session.

#	Rubric Item	ρ	p	N
4	Communicates effectively with team	-0.470	0.049	18
18	Uses patient identifiers	-0.455	0.033	22
19	Utilizes standardized practices	-0.377	0.083	22
21	Manages technology and equipment	-0.350	0.110	22
12	Prioritizes appropriately	-0.343	0.118	22

450 The per-item trends are broadly similar in magnitude
451 across items, consistent with the expectation that with $N =$
452 22 sessions, individual rubric items lack sufficient power to
453 differentiate statistically from one another. Two items reach
454 nominal significance: Item 18 (“Uses patient identifiers,”
455 $\rho = -0.455$, $p = 0.033$) and Item 4 (“Communicates ef-
456 fectively with team,” $\rho = -0.470$, $p = 0.049$). Patient
457 safety protocols may produce the strongest pattern because
458 students who excel in this expected behavior tend to per-
459 form additional wristband-checking and verification steps,
460 generating visually diverse frame sequences that are inher-
461 ently harder to classify.

462 Items related to purely procedural tasks performed in
463 a static, repetitive manner (e.g., Item 1 “Obtains pertinent
464 data,” $\rho \approx 0$) show no association, consistent with the ex-
465 pectation that these actions appear visually similar regard-
466 less of competency level. Overall, these patterns suggest
467 that vision-based analysis is most informative for expected
468 behaviors tied to procedural diversity and protocol com-
469 plexity, but is fundamentally limited in capturing the *con-*
470 *tent* of verbal communication and clinical reasoning.

5.3. RQ3: Temporal Patterns of High vs. Low Performers

471 To understand what distinguishes high- from low-
472 performing students, we partition sessions by median com-
473 petency score and construct process models from ground-
474 truth action sequences (Fig. 4).
475
476

477 Several structural differences emerge (detailed analysis
478 in App. D). Low performers show a higher Screen self-
479 loop (48% vs. 41%), reflecting more time lingering on
480 the bedside monitor, a visually uniform action that inflates
481 MOF. High performers distribute transitions more evenly
482 across Examination, Writing, and Calculator. The medi-
483 cation pathway also differs: high performers show a di-
484 rect Prep Med \rightarrow Apply Med transition (46%), while low
485 performers route through Screen (38%), suggesting work-
486 flow hesitation. High performers engage in more Exam-
487 ination actions (36 vs. 29), which involve diverse move-
488 ments that are inherently harder to classify, consistent
489 with the observed negative trend between accuracy and
490 competency. Furthermore, the low-performer model con-
491 tains more group-unique (red) transitions, indicating irreg-
492 ular workflow paths, while high performers follow a more
493 protocol-consistent progression. Finally, high performers
494 exhibit a strong Hygiene \rightarrow Screen transition (76%), sug-
495 gesting more consistent infection-control practices.

496 To rule out annotation artifacts as the source of this
497 negative trend, we perform a partial association analysis
498 controlling for six potential confounders (annotation cov-
499 erage, segment count, unique action types, average seg-
500 ment duration, video duration, and total annotations). The
501 MOF-competency association persists across all controls
502 and *strengthens* when controlling for annotation coverage
503 ($\rho: -0.439 \rightarrow -0.546$, $p = 0.009$); full results are reported
504 in App. E.

6. Discussion: Rethinking Eval Metrics

505 Our findings raise the question of whether higher frame-
506 level accuracy is always the appropriate optimization target
507 for action recognition in educational settings. In our data,
508 the classifier tends to perform better on sessions with repeti-
509 tive actions, whereas the more fluid and adaptive workflows
510 associated with higher instructor-rated competency appear
511 harder to recognize. This asymmetry is partly rooted in
512 the prototype-based design: each action class is represented
513 by a single centroid, which favors within-class visual con-
514 sistency. Students who perform an action similarly across
515 instances produce tighter feature clusters that are easier to
516 match, whereas students who vary their approach across
517 instances produce more dispersed features that weaken proto-
518 type fit. One interpretation is that competency, as assessed
519 by clinical educators, includes behavioral diversity and pro-
520 cedural flexibility that current vision models do not fully
521

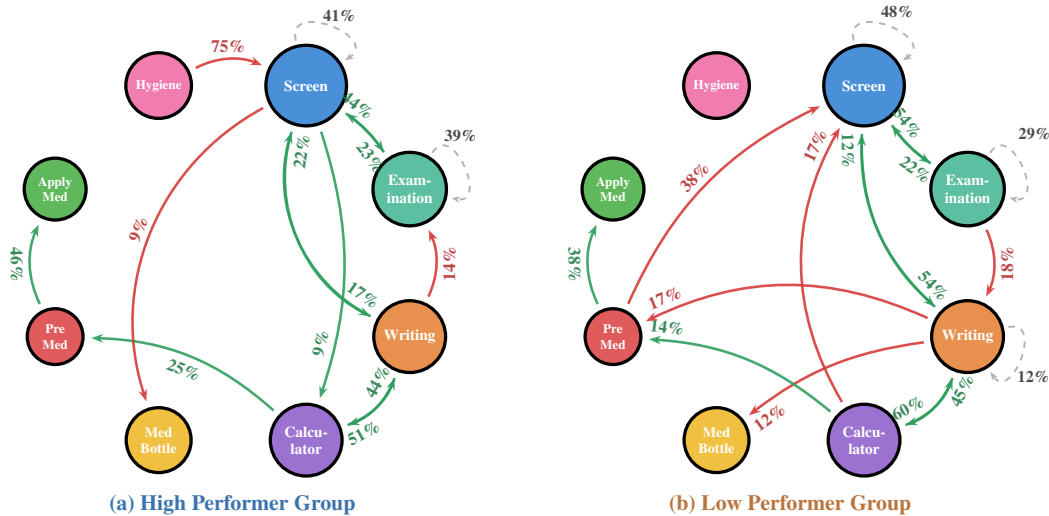


Figure 4. Process models comparing high- and low-competency groups from the ground truth of actions. The 16 fine-grained clinical actions (App. C; background excluded) are aggregated into 8 macro-categories: Examination (Palpate Wrist, Apical Pulse, Lung Sounds, Temperature, Blood Pressure), Hygiene (Hand Hygiene, Gloves), Screen (Patient History, Vital Signs), Writing, Calculator, Med Bottle, Prep Med, and Apply Med. **Green** edges denote transitions shared by both groups; **Red** edges are unique to one group. Percentages indicate transition probabilities in the arrow’s direction.

522 capture. This view also aligns with broader work suggest- 522
 523 ing that the extent to which an individual’s responses align 523
 524 with a group can relate to learning and memory outcomes; 524
 525 here, the analogous notion of “fit” is between a student’s 525
 526 action sequence and a prototype-based model constructed 526
 527 from peers through leave-one-out cross-sample prototypes. 527

528 This observation suggests a practical two-tier approach 528
 529 to automated assessment: (1) the predicted action time- 529
 530 line provides a coarse behavioral summary of what the stu- 530
 531 dent did, and (2) the recognition *difficulty* of each session, 531
 532 quantified by mIoU or F1, may serve as a complementary 532
 533 signal of holistic competency. For medication administra- 533
 534 tion, where competency is inherently sequential and correct 534
 535 actions performed in the wrong order can still constitute 535
 536 clinical error, this perspective may help educators identify 536
 537 students whose workflows deviate from the expected pro- 537
 538 cedural pathway even when checklist ratings appear sim- 538
 539 ilar. The limited variability in instructor C-CEI ratings fur- 539
 540 ther suggests that checklist-based instruments may lack the 540
 541 granularity to distinguish students with clustered overall 541
 542 scores; recognition accuracy may complement such instru- 542
 543 ments by capturing differences in *how* workflows are ex- 543
 544 ecuted. More broadly, both rubric-based assessment and 544
 545 vision-based analysis are limited to observable behavior 545
 546 and do not capture the clinical reasoning behind procedu- 546
 547 ral choices. Combining recognition accuracy with post- 547
 548 simulation reflection data, such as structured debriefs or 548
 549 self-assessments, may therefore provide a more complete 549
 550 view of student ability across behavioral and cognitive di- 550
 551 mensions. However, with only $N = 22$ sessions, this 551

522 observation remains exploratory, and recognition accuracy 522
 523 should be viewed as one potential indicator rather than a 523
 524 definitive measure. 524

525 7. Conclusion

526 We presented a three-stage framework for automated com- 526
 527 petency assessment from egocentric nursing simulation 527
 528 videos. Our results suggest that recognition accuracy may 528
 529 carry a pedagogically informative signal: more competent 529
 530 students produce diverse workflows that are systematically 530
 531 harder to classify. This suggests a two-tier assessment in 531
 532 which predicted action timelines provide a behavioral sum- 532
 533 mary and recognition difficulty provides a complementary 533
 534 competency signal, establishing the utility and boundaries 534
 535 of unimodal video-based assessment. 535

536 The primary limitation is the sample size, which reflects 536
 537 the constraints of privacy-regulated clinical data collection 537
 538 and expert annotation; larger multi-site cohorts are needed 538
 539 for generalizability. Vision alone cannot capture verbal 539
 540 communication or the clinical reasoning underlying procedu- 540
 541 ral choices. Integrating audio, gaze, and physiological 541
 542 modalities, along with post-simulation reflection data such 542
 543 as debrief transcripts, could yield a more complete picture 543
 544 of competency across behavioral, cognitive, and metacogni- 544
 545 tive dimensions. Whether the relationship between recogni- 545
 546 tion accuracy and competency extends to other educational 546
 547 domains remains an open question. 547

578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634**References**

- [1] Katie Anne Adamson. *Assessing the reliability of simulation evaluation instruments used in nursing education: A test of concept study*. Washington State University, 2011. 1
- [2] Paulo Blikstein and Marcelo Worsley. Multimodal learning analytics and education data mining: Using computational technologies to measure complex learning tasks. *Journal of Learning Analytics*, 3(2):220–238, 2016. 1
- [3] Brandon M Booth, Nigel Bosch, and Sidney K D’Mello. Engagement detection and its applications in learning: a tutorial and selective review. *Proceedings of the IEEE*, 111(10):1398–1422, 2023. 1
- [4] Nigel Bosch, Sidney D’Mello, Ryan Baker, Jaclyn Ocumpaugh, Valerie Shute, Matthew Ventura, Lubin Wang, and Weinan Zhao. Automatic detection of learning-centered affective states in the wild. In *Proceedings of the 20th international conference on intelligent user interfaces*, pages 379–388, 2015. 3
- [5] João Carreira and Andrew Zisserman. Quo vadis, action recognition? A new model and the kinetics dataset. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6299–6308, 2017. 2
- [6] Clayton Cohn, Caitlin Snyder, Joyce Horn Fonteles, Ashwin TS, Justin Montenegro, and Gautam Biswas. A multimodal approach to support teacher, researcher and ai collaboration in stem+ c learning environments. *British Journal of Educational Technology*, 56(2):595–620, 2025. 1
- [7] Eduardo Davalos, Yike Zhang, Ashwin T S, Joyce Horn Fonteles, Umesh Timalisina, and Gautam Biswas. 3d gaze tracking for studying collaborative interactions in mixed-reality environments. In *Companion Proceedings of the 26th International Conference on Multimodal Interaction*, pages 175–183, 2024. 1
- [8] Guodong Ding, Fadime Sener, and Angela Yao. Temporal action segmentation: An analysis of modern techniques. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(2):1011–1030, 2024. 2, 3
- [9] Sidney K D’mello and Jacqueline Kory. A review and meta-analysis of multimodal affect detection systems. *ACM computing surveys (CSUR)*, 47(3):1–36, 2015. 1
- [10] K Anders Ericsson, Ralf Th Krampe, and Clemens Tesch-Römer. The role of deliberate practice in the acquisition of expert performance. *Psychological Review*, 100(3):363–406, 1993. 1
- [11] Joyce Fonteles, Eduardo Davalos, TS Ashwin, Yike Zhang, Mengxi Zhou, Efrat Ayalon, Alicia Lane, Selena Steinberg, Gabriella Anton, Joshua Danish, et al. A first step in using machine learning methods to enhance interaction analysis for embodied learning environments. In *International Conference on Artificial Intelligence in Education*, pages 3–16. Springer, 2024. 3
- [12] Joyce Horn Fonteles, Clayton Cohn, Efrat Ayalon, Mengxi Zhou, Ashwin TS, Eduardo Davalos, Zhijian Li, Surya Rayala, Divya Mereddy, Austin Coursey, et al. Analyzing embodied learning in classroom settings: A human-in-the-loop ai approach for multimodal learning analytics. *Learning and Instruction*, 103:102274, 2026. 3
- [13] Isabel Funke, Sjoerd T Mees, Jürgen Weitz, and Stefanie Speidel. Video-based surgical skill assessment using 3D convolutional neural networks. *International Journal of Computer Assisted Radiology and Surgery*, 14(7):1217–1225, 2019. 3
- [14] Emmett D Goodman, Krishna K Patel, Yilun Zhang, William Locke, Chris J Kennedy, Rohan Mehrotra, Stephen Ren, Melody Y Guan, Maren Downing, Hao Wei Chen, et al. A real-time spatiotemporal ai model analyzes skill in open surgical videos. *arXiv preprint arXiv:2112.07219*, 2021. 3
- [15] Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al. Ego4D: Around the world in 3,000 hours of egocentric video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18995–19012, 2022. 1, 3
- [16] Daniel A. Hashimoto, Guy Rosman, Daniela Rus, and Ozanan R. Meireles. Artificial intelligence in surgery: Promises and perils. *Annals of Surgery*, 268(1):70–76, 2018. 3
- [17] Jennifer K Hayden, Richard A Smiley, Maryann Alexander, Suzan Kardong-Edgren, and Pamela R Jeffries. The NCSBN national simulation study: A longitudinal, randomized, controlled study replacing clinical hours with simulation in pre-licensure nursing education. *Journal of Nursing Regulation*, 5(2):C1–S64, 2014. 1
- [18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. 4
- [19] Sanchit Hira, Digvijay Singh, Tae Soo Kim, Shobhit Gupta, Gregory Hager, Shameema Sikder, and S Swaroop Vedula. Video-based assessment of intraoperative surgical skill. *International journal of computer assisted radiology and surgery*, 17(10):1801–1811, 2022. 6
- [20] Sissel Eikeland Husebø, Febe Friberg, Eldar Søreide, and Hans Rystedt. Instructional problems in briefings: How to prepare nursing students for simulation-based cardiopulmonary resuscitation training. *Clinical Simulation in Nursing*, 9(8):e307–e318, 2013. 1
- [21] Stephen Hutt, Kristina Krasich, Caitlin Mills, Nigel Bosch, Shelby White, James R Brockmole, and Sidney K D’Mello. Automated gaze-based mind wandering detection during computerized learning in classrooms: S. hutt et al. *User Modeling and User-Adapted Interaction*, 29(4):821–867, 2019. 1
- [22] Pamela R Jeffries. A framework for designing, implementing, and evaluating simulations used as teaching strategies in nursing. *Nursing Education Perspectives*, 26(2):96–103, 2005. 1, 2
- [23] Suzan Kardong-Edgren, Kathleen A. Adamson, and Cynthia Fitzgerald. A review of currently published evaluation instruments for human patient simulation. *Clinical Simulation in Nursing*, 6(1):e25–e35, 2010. 1
- [24] Ahmad Khalifa, Owais Tahhan, Mohammed Albazooni, Mohammed Saeed, Ruha Hamdi, Megan Stanners, Amman Malik, and Adnan Malik. Automated and artificial intelligence

- (ai)-derived performance assessment in surgical simulation: A systematic review. *Cureus*, 17(12), 2025. 2
- [25] J. Richard Landis and Gary G. Koch. The measurement of observer agreement for categorical data. *Biometrics*, 33(1): 159–174, 1977. 3
- [26] Kathie Lasater. Clinical judgment development: Using simulation to create an assessment rubric. *Journal of Nursing Education*, 46(11):496–503, 2007. 1
- [27] Michele Mazzamuto, Antonino Furnari, Yoichi Sato, and Giovanni Maria Farinella. Gazing into missteps: Leveraging eye-gaze for unsupervised mistake detection in egocentric videos of skilled human activities. In *European Conference on Computer Vision (ECCV)*, 2024. 1
- [28] Hans Moravec. *Mind Children: The Future of Robot and Human Intelligence*. Harvard University Press, Cambridge, MA, 1988. 6
- [29] Chinedu Innocent Nwoye, Tong Yu, Cristians Gonzalez, Barbara Seez, Didier Mutter, Jacques Marescaux, and Nicolas Padoy. Rendezvous: Attention mechanisms for the recognition of surgical action triplets in endoscopic videos. page 102433, 2022. 3
- [30] Xavier Ochoa and Marcelo Worsley. Augmenting learning analytics with multimodal sensory data. *Journal of Learning Analytics*, 3(2):213–219, 2016. 3
- [31] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafranec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. DINOv2: Learning robust visual features without supervision. *Transactions on Machine Learning Research*, 2024. 3, 4
- [32] James W. Pellegrino, Naomi Chudowsky, and Robert Glaser. *Knowing What Students Know: The Science and Design of Educational Assessment*. National Academies Press, 2001. 1
- [33] Lawrence R Rabiner. A tutorial on hidden Markov models and selected applications in speech recognition. *Proceedings of the IEEE*, 77(2):257–286, 1989. 3, 4
- [34] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *Proceedings of the 38th International Conference on Machine Learning (ICML)*, pages 8748–8763, 2021. 4
- [35] Rajiv Ranganathan, Mei-Hua Lee, and Karl M Newell. Repetition without repetition: Challenges in understanding behavioral flexibility in motor skill. *Frontiers in Psychology*, 11:2018, 2020. 6
- [36] Ginger Schroers and Jill Pfeiffer. Tool development and testing: An objective measurement of medication administration competency. *Nursing Education Perspectives*, 46(5):E37–E39, 2025. 2, 1
- [37] Deborah Schwengel, Ignacio Villagrán, Geoffrey Miller, Constanza Miranda, et al. Multimodal assessment in clinical simulations: A guide for moving towards precision education. *Medical Science Educator*, 35(2):1025–1034, 2024. 1
- [38] Saswat Shankar, Caitlyn Ruiz, Zijian Zheng, Marcelo Worsley, and Paulo Blikstein. Multimodal methods for analyzing learning and training environments: A systematic literature review. *arXiv preprint arXiv:2408.14491*, 2024. 3
- [39] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017. 3, 4
- [40] Caitlin Snyder, Nicole M Hutchins, Clayton Cohn, Joyce Horn Fonteles, and Gautam Biswas. Analyzing students collaborative problem-solving behaviors in synergistic stem+ c learning. In *Proceedings of the 14th learning analytics and knowledge conference*, pages 540–550, 2024. 3
- [41] Abed Soleymani, Ali Akbar Sadat Asl, Mojtaba Yeganejou, Scott Dick, Mahdi Tavakoli, and Xingyu Li. Surgical skill evaluation from robot-assisted surgery recordings. In *2021 International Symposium on Medical Robotics (ISMR)*, pages 1–6. IEEE, 2021. 2, 6
- [42] Ekta Sood, Fabian Kögel, Philipp Müller, Dominike Thomas, Mihai Băce, and Andreas Bulling. Multimodal integration of human-like attention in visual question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2648–2658, 2023. 3
- [43] Ömer Sümer, Patricia Goldberg, Sidney D’Mello, Peter Gerjets, Ulrich Trautwein, and Enkelejda Kasneci. Multimodal engagement analysis from facial videos in the classroom. *IEEE Transactions on Affective Computing*, 14(2):1012–1027, 2021. 1
- [44] Michael Todd, Julie A. Manz, Kathleen S. Hawkins, Michele E. Parsons, and Mary Hercinger. The development of a quantitative evaluation tool for simulations in nursing education. *International Journal of Nursing Education Scholarship*, 5(1), 2008. 1
- [45] Andru Putra Twinanda, Sherif Shehata, Didier Mutter, Jacques Marescaux, Michel de Mathelin, and Nicolas Padoy. EndoNet: A deep architecture for recognition tasks on laparoscopic videos. *IEEE Transactions on Medical Imaging*, 36(1):86–97, 2017. 3
- [46] Caleb Vatral, Gautam Biswas, and Benjamin Goldberg. A theoretical framework for performance analysis in competency-based experiential learning environments. In *AI and Gamification Technologies for Complex Work*, pages 1–19. CRC Press, 2025. 1
- [47] Hung-Hsuan Yen, Ming-Chih Ho, Yi-Hsiang Hsiao, and Chun-Chieh Huang. Surgical video-based temporal action analysis algorithm and competency assessment in laparoscopic cholecystectomy: development and exploratory evaluation. *Surgical Endoscopy*, 2025. 2

AI-Assisted Competency Assessment from Egocentric Video in Simulation-Based Nursing Education

Supplementary Material

797	A. Simulation Scenario Summary		
798	Scenario Context and Setting	The simulation scenario	
799		and debrief were created by a nursing teaching instructor and have been used in nursing school classroom settings. The simulation scenario takes place in a high-fidelity pediatric emergency room bay. A standardized pediatric manikin representing a toddler and a faculty facilitator acting as the patient's caregiver are present at the bedside. The scenario is designed to evaluate pediatric assessment, weight-based medication administration, and caregiver communication competencies.	
800			
801			
802			
803			
804			
805			
806			
807			
808	Anonymized Patient Profile	The simulated patient is a 16-month-old toddler (9.6 kg, 76 cm) presenting with a primary diagnosis of croup (laryngotracheobronchitis). The caregiver reports a 3-day history of upper respiratory infection symptoms, with sudden overnight onset of a barking cough, hoarse voice, and inspiratory stridor. On arrival, the patient is placed on continuous pulse oximetry, heart rate, and respiratory rate monitoring, and maintained on humidified oxygen at 1 LPM via pediatric face mask.	
809			
810			
811			
812			
813			
814			
815			
816			
817	Simulation Learning Objectives	The scenario targets four core nursing competencies: (1) performing a focused pediatric assessment while maintaining age-appropriate patient safety; (2) recognizing pediatric fever and calculating accurate weight-based dosages for oral antipyretic medications; (3) preparing and administering pediatric oral suspensions safely; and (4) providing clear, developmentally appropriate education to caregivers regarding at-home medication administration.	
818			
819			
820			
821			
822			
823			
824			
825			
826	Scenario Progression and Key Interventions	The simulation unfolds across three phases. In Phase 1 , the student initiates care by performing hand hygiene, verifying patient identification using two identifiers, and introducing themselves to the caregiver. Initial vitals reflect tachypnea and tachycardia consistent with the patient's respiratory distress. In Phase 2 , the student performs a focused respiratory assessment, noting mild expiratory wheezing and intermittent barking cough. A bedside temperature check reveals a fever of 102.6°F, prompting the student to review physician orders and perform a weight-based medication calculation	
827			
828			
829			
830			
831			
832			
833			
834			
835			
836			
		for oral acetaminophen suspension (160 mg/5 mL):	837
		$15 \text{ mg/kg} \times 9.6 \text{ kg} = 144 \text{ mg},$	838
		$144 \text{ mg} \times \frac{5 \text{ mL}}{160 \text{ mg}} = 4.5 \text{ mL}.$	(3) 839
	In Phase 3 , after preparing the medication in an amber oral dosing syringe, the student engages in targeted caregiver education. Key instructional points include advising against using household spoons for measuring, demonstrating correct syringe administration technique to prevent choking, and establishing safe guidelines for dosing frequency at home.		840
			841
			842
			843
			844
			845
			846
	B. Instructor Competency Rubric		847
	The instructor rubric is an adapted version of the Creighton Competency Evaluation Instrument (C-CEI) [44], which maps 23 expected behaviors to broader concepts of competency (e.g., clinical judgment, patient safety, communication). Each item is rated on a 1–5 scale (Poor to Exceptional). Because our study uses egocentric video without audio, only the 11 video-observable items (highlighted in Tab. 5) contribute to each student's competency percentage. The remaining 12 items require verbal or cognitive assessment not accessible from visual data alone.		848
			849
			850
			851
			852
			853
			854
			855
			856
			857
	Rationale for the video-observable subset. Items requiring verbal content (e.g., "Communicates effectively with team," "Provides evidence-based rationale") or internal cognitive processes (e.g., "Reflects on clinical experience") cannot be assessed from silent egocentric video. The 11 retained items correspond to physical actions and procedural behaviors that produce visible evidence in the video stream: checking wristbands, performing hand hygiene, documenting on screens, measuring vital signs, administering medications, and following safety protocols. This principled subset ensures that the competency score reflects only expected behaviors that our vision-based system could plausibly detect. Note that in the per-item association analysis (Tab. 4), we report associations for all 23 items to explore whether vision-based features carry any indirect signal for non-observable behaviors.		858
			859
			860
			861
			862
			863
			864
			865
			866
			867
			868
			869
			870
			871
			872
			873
	C. Action Annotation Rubric		874
	The following is from our annotation codebook, used by trained coders to produce ground-truth action annotations, and is inspired by [36]. Actions represent discrete, observ-		875
			876
			877

Table 5. Full 23-item C-CEI rubric. Each item specifies an expected behavior; highlighted rows (✓) are the 11 video-observable items used for competency scoring.

#	Item Description	Video?
1	Obtains pertinent data	✓
2	Performs follow-up assessments as needed	✓
3	Assesses the environment	×
4	Communicates effectively with team	×
5	Communicates effectively with patient	×
6	Documents clearly, concisely, and accurately	✓
7	Responds to abnormal findings appropriately	×
8	Promotes professionalism	×
9	Interprets vital signs	✓
10	Interprets laboratory results	×
11	Interprets subjective/objective data	×
12	Prioritizes appropriately	✓
13	Performs evidence-based interventions	✓
14	Provides evidence-based rationale for interventions	×
15	Evaluates evidence-based interventions and outcomes	×
16	Reflects on clinical experience	×
17	Delegates appropriately	×
18	Uses patient identifiers	✓
19	Utilizes standardized practices and precautions	✓
20	Administers medications safely	✓
21	Manages technology and equipment	✓
22	Performs procedures correctly	✓
23	Reflects on potential hazards and errors	×

878 able physical behaviors; verbal introductions are captured
879 separately by the Communication layer.

880 General coding rules.

- 881 1. Code only what is directly observable; do not infer in-
882 tent.
- 883 2. When in doubt, leave the segment unlabeled.
- 884 3. Annotations must not overlap within the Action layer.
- 885 4. Start when the action begins (first observable move-
886 ment); end when it concludes (hands leave the object,
887 body repositions away).
- 888 5. Brief interruptions (<2 s): code as one continuous seg-
889 ment.

890 **Action definitions.** Tab. 6 lists the $K=16$ fine-grained clin-
891 ical action classes. Frames that do not correspond to any
892 of these classes (e.g., walking, adjusting equipment, idle
893 periods between clinical actions) are left unannotated and
894 treated as the background class a_{\emptyset} , yielding $K+1=17$ la-
895 bels in total for recognition.

896 **Disambiguation guidelines.** Several action pairs are visu-
897 ally similar and require explicit decision rules:

898 *Lung Sounds (#8) vs. Apical Pulse (#9):* Stethoscope on the
899 back or moved across multiple chest positions is coded as
900 #8. Stethoscope held at the left chest apex in one position
901 for ≥ 15 s is coded as #9. If placement is unclear, default to
902 #8 and flag for review.

903 *Calculator (#13) vs. Phone (#14):* Tapping numbers on a

Table 6. The $K=16$ clinical action classes used for temporal anno-
tation and few-shot recognition, with brief operational definitions.
An additional background class a_{\emptyset} (not shown) captures all non-
clinical frames, yielding 17 labels total.

ID	Action Class	Definition
1	Perform Hand Hygiene	Uses hand sanitizer or washes hands at sink
2	Put on Gloves	Retrieves and dons disposable gloves
3	Check Patient Wristband	Visually inspects or scans patient wristband
4	Check Patient History Screen	Reads electronic health record on screen
5	Examine Med Bottle	Picks up and reads medication label
6	Review Vital Signs Screen	Reads the vital signs monitor (HR, BP, SpO ₂)
7	Assess Vital Signs (Palpate Wrist)	Manually palpates radial pulse
8	Auscultate Lung Sounds	Places stethoscope on chest/back for breath sounds
9	Measure Apical Pulse	Places stethoscope at heart apex, held ≥ 15 s
10	Measure Temperature	Uses thermometer (oral, tympanic, temporal)
11	Measure Blood Pressure	Initiates BP reading via monitor or manual cuff
12	Writing	Pen-to-paper: notes, calculations, forms
13	Use Calculator	Computes dosage on physical or phone calculator
14	Check Phone	Interacts with phone for non-calculator purposes
15	Prepare Medication	Draws syringe, crushes tablet, mixes solution
16	Apply Medication to Patient	Administers medication: oral, IV, injection, topical

calculator app or physical calculator is #13. Scrolling, read- 904
ing, or swiping on a phone (non-calculator) is #14. 905

Patient History Screen (#4) vs. Vital Signs Screen (#6): If 906
the screen shows waveforms or real-time numeric readings, 907
code as #6. If it shows text-based records, history, or med- 908
ication orders, code as #4. Pressing a button on the vitals 909
monitor to initiate a BP measurement is coded as #11. 910

911 D. Process Model Details

The five key structural differences between high- and low- 912
performer process models (Fig. 4) are elaborated below. 913

Screen self-loop. Low-performing students exhibit a higher 914
self-loop on the Screen action (48% vs. 41%), spending 915
proportionally more time returning to the bedside moni- 916
tor without transitioning to other clinical actions. High 917
performers distribute transitions away from Screen more 918
evenly across Examination, Writing, and Calculator, reflect- 919
ing a more fluid workflow. Screen actions are visually static 920
and uniform, making them easy for the classifier and inflat- 921
ing MOF for the LOW group. 922

Medication pathway. High performers show a strong di- 923
rect Prep Med \rightarrow Apply Med transition (46%), indicating a 924
coherent prepare-then-administer sequence. Low perform- 925
ers lack this link; instead, Prep Med routes back to Screen 926
(38%), suggesting hesitation or uncertainty in the medica- 927
tion procedure. 928

Examination frequency. High performers engage in more 929
Examination actions (36 vs. 29), while low performers pro- 930
duce more Writing and Screen actions (42 and 79 vs. 37 and 931
74). Physical examination (lung sounds, blood pressure, 932
palpation) involves diverse movements inherently harder 933
to classify, consistent with the observed negative trend be- 934
tween accuracy and competency. 935

Transition irregularity. The low-performer model con- 936
tains more group-unique (red) transitions, indicating ir- 937

938 regular workflow paths. High performers follow a more
939 protocol-consistent progression with fewer idiosyncratic
940 transitions.

941 **Hygiene compliance.** Hygiene actions connect to Screen
942 with 76% probability in high performers, suggesting consis-
943 tent hand hygiene before engaging with the patient monitor.
944 This transition is less prominent in low performers, pointing
945 to less consistent infection control practices.

946 These process model comparisons offer actionable in-
947 sight for clinical educators: the transition graphs visualize
948 where each student’s workflow diverges from the expected
949 clinical pathway, enabling targeted remediation of specific
950 procedural gaps.

951 E. Annotation Confound Analysis

952 A potential concern is that annotation artifacts drive the neg-
953 ative trend between classification accuracy and competency,
954 since HIGH competency sessions have lower annotation
955 coverage (40% vs. 50%). We perform partial association
956 analysis, controlling for six potential confounds (Tab. 7).
957 If any drove the observed pattern, controlling for it would
958 weaken or eliminate the effect.

Table 7. Robustness analysis. *Partial ρ* : Spearman association between MOF and competency after controlling for each variable. *Var \leftrightarrow MOF ρ* : bivariate association between each variable and MOF. The MOF–competency association persists across all controls and *strengthens* when controlling for annotation coverage (bolded). No control variable independently predicts MOF (all $p > 0.46$).

Control Variable	Partial ρ	p	Var \leftrightarrow MOF ρ	p
None (baseline)	−0.439	0.041	—	—
Annotation coverage	−0.546	0.009	−0.032	0.887
# GT action segments	−0.427	0.047	+0.115	0.611
# Unique GT action types	−0.438	0.041	+0.027	0.907
Avg segment duration	−0.437	0.042	−0.165	0.462
Video duration	−0.454	0.034	+0.074	0.744
# All annotations	−0.427	0.047	+0.115	0.611

959 The pattern persists across all controls. When con-
960 trolling for annotation coverage, the effect *strengthens* (ρ :
961 $−0.439 \rightarrow −0.546$), and no control variable independently
962 predicts MOF (all $p > 0.46$), confirming that the negative
963 trend reflects workflow complexity rather than annotation
964 density.

965 F. Inter-Rater Reliability

966 A second rater independently annotated 3 stratified videos
967 (low / median / high competency) to assess annotation reli-
968 ability. Agreement was measured using frame-level Cohen’s
969 κ at 1 Hz resolution. To avoid inflation from unannotated
970 frames, κ was computed only over frames where at least
971 one rater placed a label.

Mean $\kappa = 0.708 \pm 0.199$ (substantial agreement; [25]).
As a secondary metric, mean per-class IoU = 0.697 ± 0.143 ,
and both raters identified identical action type sets in all
3 videos (Jaccard = 1.0). Disagreements were predomi-
nantly in segment boundary placement, particularly action
endpoints (mean $|\Delta| = 3.4$ s), rather than action identifica-
tion or ordering. This pattern indicates that raters agree on
which actions occur and *in what order*, with variability con-
fined to the precise temporal boundaries, consistent with the
known difficulty of endpoint annotation in temporal action
segmentation [8].

972
973
974
975
976
977
978
979
980
981
982