

Finanziato dall'Unione europea NextGenerationEU













Spin-off of the University of Catania



Procedure Understanding from Egocentric Videos

Giovanni Maria Farinella

FPV @ Image Processing Laboratory - <u>http://iplab.dmi.unict.it/fpv</u>

Next Vision - <u>http://www.nextvisionlab.it/</u>

Department of Mathematics and Computer Science - University of Catania

giovanni.farinella@unict.it - www.dmi.unict.it/farinella



Egocentric vision or first-person vision is a sub-field of computer vision that entails analyzing images, videos captured by a wearable device, which is typically worn on the head or on the chest and naturally approximates the visual field of the camera wearer.

Consequently, visual data capture the part of the scene on which the user focuses to carry out the task at hand and offer a valuable perspective to understand the user's intents and activities and their context in a naturalistic setting.

This research area is sometimes referred with the name "Wearable Vision".

Università What's Relevant in EgoVision? A top-down approach

An Outlook into the Future of Egocentric Vision



Abstract What will the future be? We wonder! In this survey, we explore the gap between current research in egocentric vision and the ever-anticipated future, where wearable computing, with outward facing cameras and digital overlays, is expected to be integrated in our every day lives. To understand this gap, the article starts by envisaging the future through character-based stories, showcasing through examples the limitations of current technology. We then provide a mapping between this future and previously defined research tasks. For each task, we survey its seminal works, current stateof-the-art methodologies and available datasets, then reflect on shortcomings that limit its applicability to future research. Note that this survey focuses on software models for egocentric vision, independent of any specific hardware. The paper concludes with recommendations for areas of immediate explorations so as to unlock our path to the future always-on, personalised and life-enhancing egocentric vision.

Keywords Egocentric Vision, Future, Survey, Localisation, Scene Understanding, Recognition, Anticipation, Gaze Prediction, Social Understanding, Body Pose Estimation, Hand and Hand-Object Interaction, Person Identification, Summarisation, Dialogue, Privacy

Contents

- *: Equal Contribution/First Author
- [†]: Equal Senior Author

C. Plizzari, G. Goletto and T. Tommasi, Politecnico di Torino, Italy · A. Furnari, F. Ragusa and G. M. Farinella, University of Catania, Italy · S. Bansal and D. Damen, University of Bristol, UK. E-mail: Tatiana.Tommasi@polito.it

	2.1 EGO-Home	
	2.2 EGO-Worker	
	2.3 EGO-Tourist	
	2.4 EGO-Police	
	2.5 EGO-Designer	
3	From Narratives to Research Tasks 8	
4	Research Tasks and Capabilities	
	4.1 Localisation	
	4.2 3D Scene Understanding	
	4.3 Recognition	
	4.4 Anticipation	
	4.5 Gaze Understanding and Prediction 23	
	4.6 Social Behaviour Understanding	
	4.7 Full-body Pose Estimation	
	4.8 Hand and Hand-Object Interactions 30	
	4.9 Person Identification	
	4.10 Summarisation	
	4.11 Dialogue	
	4.12 Privacy	
	4.13 Beyond individual tasks	
5	General Datasets	
6	Conclusion 49	

1 Introduction

Designing and building tools able to support human activities, improve quality of life, and enhance individuals' abilities to achieve their goals is the ever-lasting aspiration of our species. Among all inventions, digital computing has already had a revolutionary effect on human history. Of particular note is mobile technology, currently integrated in our lives through hand-held devices, i.e. *mobile smart phones*. These are nowadays the de facto for outdoor navigation, capturing static and moving footage of our everyday and connecting us to both familiar and novel connections and experiences.

However, humans have been dreaming about the next-version of such mobile technology — wearable computing, for a considerable amount of time. Imaginations



A lot of data!



Rather than being extensive, we considered **seminal** and **state-of-the-art** works

Plizzari, C., Goletto, G., Furnari, A., Bansal, S., Ragusa, F., Farinella, G. M., Damen., D. & Tommasi, T. (2023). An Outlook into the Future of Egocentric Vision. International Journal of Computer Vision.

Università di Catania What's Relevant in Egovision? Scenarios

An Outlook into the Future of Egocentric Vision



Abstract What will the future be? We wonder! In this survey, we explore the gap between current research in egocentric vision and the ever-anticipated future, where wearable computing, with outward facing cameras and digital overlays, is expected to be integrated in our every day lives. To understand this gap, the article starts by envisaging the future through character-based stories, showcasing through examples the limitations of current technology. We then provide a mapping between this future and previously defined research tasks. For each task, we survey its seminal works, current stateof-the-art methodologies and available datasets, then reflect on shortcomings that limit its applicability to future research. Note that this survey focuses on software models for egocentric vision, independent of any specific hardware. The paper concludes with recommendations for areas of immediate explorations so as to unlock our path to the future always-on, personalised and life-enhancing egocentric vision.

Keywords Egocentric Vision, Future, Survey, Localisation, Scene Understanding, Recognition, Anticipation, Gaze Prediction, Social Understanding, Body Pose Estimation, Hand and Hand-Object Interaction, Person Identification, Summarisation, Dialogue, Privacy

Contents

 1
 Introduction
 1

 2
 Imagining the Future
 2

*: Equal Contribution/First Author

[†]: Equal Senior Author

C. Plizzari, G. Goletto and T. Tommasi, Politecnico di Torino, Italy · A. Furnari, F. Ragusa and G. M. Farinella, University of Catania, Italy · S. Bansal and D. Damen, University of Bristol, UK. E-mail: Tatiana.Tommasi@polito.it

2.2 2.5 EGO-Designer 4.5 Gaze Understanding and Prediction 23 4.8 Hand and Hand-Object Interactions 30 Conclusion 49

1 Introduction

Designing and building tools able to support human activities, improve quality of life, and enhance individuals' abilities to achieve their goals is the ever-lasting aspiration of our species. Among all inventions, digital computing has already had a revolutionary effect on human history. Of particular note is mobile technology, currently integrated in our lives through hand-held devices, i.e. *mobile smart phones*. These are nowadays the de facto for outdoor navigation, capturing static and moving footage of our everyday and connecting us to both familiar and novel connections and experiences.

However, humans have been dreaming about the next-version of such mobile technology — wearable computing, for a considerable amount of time. Imaginations





Plizzari, C., Goletto, G., Furnari, A., Bansal, S., Ragusa, F., Farinella, G. M., Damen., D. & Tommasi, T. (2023). An Outlook into the Future of Egocentric Vision. International Journal of Computer Vision.

Università **Procedure Understanding**

Procedural activities are sequences of key-steps aimed at achieving specific goals.





Task: Given a video segment s_t and its previous video segment history, models have to: 1) determine previous keysteps (to be performed before s_t); infer if s_t is 2) optional or 3) a procedural mistake; 4) predict missing keysteps (should have been performed before s_t but were not); and 5) next keysteps (for which dependencies are satisfied).

^M Università Lots of Ego-data available for Procedure Understanding



MECCANO 💐



EgoProceL



IndustReal



EgoHOS



EPIC-Tent









EPIC-KITCHENS VISOR



ENIGMA-51





HoloAssist



Assembly101

Università di Catania Among Downstream Tasks in Procedure Understanding





Università di Catania Among the Tasks in Procedure Understanding



Human-Object Interaction Detection



Online Mistake Detection







NEXT VISI ⁽)/₍N



Spin-off of the University of Catania

Are Synthetic Data Useful for Egocentric Hand-Object Interaction Detection?





Poster Session 2 4:30PM TUE 1 OCT

Data and Code: https://fpv-iplab.github.io/HOI-Synth/

Leonardi, R., Furnari, A., Ragusa, F., & Farinella, G. M. (2024). Are Synthetic Data Useful for Egocentric Hand-Object Interaction Detection? An Investigation and the HOI-Synth Domain Adaptation Benchmark. European Conference on Computer Vision (ECCV)

Università Data Generation Pipeline



HOI-Synth Dataset – Automatically Labeled



Data and Code: <u>https://fpv-iplab.github.io/HOI-Synth/</u>



Università Can synthetic data help to train models?



Data and Code: <u>https://fpv-iplab.github.io/HOI-Synth/</u>

Università di Catania Results - EPIK-Kitchens VISOR

	a) Unsupervised Setting												
% Real Labeled Data	Approach	Overall	Η	H+S	H+C	0							
0%	Synthetic-Only	09.88	28.41	24.89	08.64	01.23							
070	UDA	33.33	80.16	65.98	33.47	8.35							
Absolute Impro	vement	+23.45	+51.75	+41.09	+24.83	+7.12							
b) Semi-supervised Setting													
% Real Labeled Data	Approach	Overall	Η	H+S	H+C	Ο							
10%	Real-Only	38.55	87.45	83.27	51.98	19.47							
(3.286 images)	Synthetic+Real	37.62	86.39	82.85	52.25	<u>23.03</u>							
(3,200 mages)	SSDA	<u>44.22</u>	<u>89.05</u>	80.77	46.83	20.41							
Absolute Impro	vement	+5.67	+1.60	-0.42	+0.27	+3.56							
25%	Real-Only	37.90	90.14	85.66	53.99	17.85							
(8.215 images)	Synthetic+Real	38.19	89.98	84.67	55.88	18.49							
(8,215 mages)	SSDA	$\begin{array}{c c c c c c c c c c c c c c c c c c c $	$\underline{22.15}$										
Absolute Impro	vement	+7.65	+0.23	-0.99	+1.89	+4.30							
50%	Real-Only	38.15	91.16	86.05	52.28	17.92							
(16.420 imagos)	Synthetic+Real	43.52	91.34	85.85	54.09	19.06							
(10,429 mages)	SSDA	46.47	90.94	85.73	$\underline{58.02}$	$\underline{23.49}$							
Absolute Impro	vement	+8.32	+0.18	-0.20	+5.74	+5.57							

c) Fully-supervised Setting

	/ • -		<u> </u>			
% Real Labeled Data	Approach	Overall	Η	H+S	H+C	0
100%	Real-Only	45.33	92.25	88.54	59.24	24.23
(32.857 imp gos)	$\operatorname{Synthetic+Real}$	44.52	91.45	<u>88.94</u>	56.55	$\underline{27.77}$
(32,837 mages)	FSDA	<u>46.48</u>	91.83	87.65	57.63	24.03
Absolute Impro	+1.15	-0.42	+0.40	-1.61	+3.54	

Similar behaviour is observed for the EgoHOS and ENIGMA-51 datasets (see the paper)



Università di Catania Among the Tasks in Procedure Understanding



3







Egocentric Action Scene Graphs for Long-Form Understanding of Egocentric Videos





CODE AND DATA: https://github.com/fpv-iplab/EASG



Università di Catania **Problem**

How to represent egocentric videos for long-term understanding?



- Understand sequences of activities performed by the camera wearer in different physical locations
- Egocentric video is by its own nature long-form
- Egocentric vision systems require algorithms able to represent and process video over temporal spans that last in the order of minutes or hours
- Examples of applications are action anticipation, video summarization, and episodic memory retrieval
- Lack of a comprehensive and long-form representation of videos that algorithms can rely on
- Popular highlevel human-gathered representations being in the form of textual narrations, verb-noun action labels, temporal bounds for action segments, object bounding boxes, object state changes, and hand-object interaction states, are all short-range representations describing temporal spans lasting few seconds.

Università di Catania Egocentric Action Scene Graphs (EASG)



Extending Ego4D with Egocentric Action Scene Graph Representations



Dataset	Dynamic	Egocentric	Sequences	Hours	Avg. Len. (seconds)	Avg. Graphs per Vid.	Obj Cls	Verb Cls	Rel Cls
VidVRD [42]	×	X	1,000	3	11	3.9*	35	25**	132
VidOR [43]	×	×	10,000	99	35	8.8^* action + 29.2* spatial	80	42	50
Action Genome [49]	\checkmark	×	10,000	82	30	5	35	-	25
PVSG [51]	×	Partly (28%)	400	9	77	382	126	44	57
HOMAGE [38]	×	paired ego-exo	1,752	25	3	3.8	86	453	29
Ego4D-EASG (Ours)	\checkmark	\checkmark	221	11.4	186	28.3	407	219	16
	•		•						

CODE AND DATA: <u>https://github.com/fpv-iplab/EASG</u>

		With Constraint									No Constraint							
Method	Edge Cls				SG Cls EASG Cls		ls	Edge Cls			SG Cls			EASG Cls				
	R@10	R@20	R@50	R@10	R@20	R@50	R@10	R@20	R@50	R@10	R@20	R@50	R@10	R@20	R@50	R@10	R@20	R@50
Random Guess	8.0	8.0	8.0	0.2	0.4	1.0	0.0	0.0	0.0	36.5	72.6	99.9	0.3	0.5	1.0	0.0	0.0	0.0
Baseline (Ours)	60.4	60.4	60.4	41.4	44.3	50.6	14.3	16.4	17.9	94.4	99.8	100	51.6	58.2	62.4	14.7	18.3	20.9

Verb-features: extracted with SlowFast model

Bbox-features: Faster-RCNN RoIAlign features





Università di Catania EASG Generation - Qualitative Result Examples









Action Anticipation

Summarization

			Ve	erb	No	oun	Action		
	Seq. length T	Avg. duration	Top-1	Top-5	Top-1	Top-5	Top-1	Top-5	
V-N	5	19s	2.54	5.01	47.68	62.24	1.28	2.60	
EASG	5	19s	3.33	<u>9.53</u>	48.84	<u>66.03</u>	1.88	5.24	
V-N	20	82s	<u>3.43</u>	8.41	46.69	64.85	2.01	4.98	
EASG	20	82s	5.94	15.97	47.36	67.26	3.40	9.24	
Improv	vement		+2.51	+7.56	+0.67	+2.41	+1.39	+4.26	

	CIDEr	ROUGE-1	METEOR
V-N	9.42	31.5	26.09
EASG	13.79	33.3	26.30
Narrations	19.99	37.7	29.43



CODE AND DATA: <u>https://github.com/fpv-iplab/EASG</u>

Università di Catania Among the Tasks in Procedure Understanding









PREGO: online mistake detection in PRocedural EGOcentric videos



Alessandro Flaborea



Guido Maria D'Amely di Leonardo Plini Melendugno, Ph.D.

Luca Scofano



Edoardo De Matteis



Antonino Furnari



Fabio Galasso



Alessandro Flaborea, Guido D'Amely, Leonardo Plini, Luca Scofano, Edoardo De Matteis, Antonino Furnari, Giovanni Maria Farinella, Fabio Galasso (2024). PREGO: online mistake detection in PRocedural EGOcentric videos. In Conference on Computer Vision and Pattern Recognition (CVPR 2024)

- Understand when the user makes a mistake to assist them;
- Current methods rely on labeled mistakes, which are hard to obtain;
- We argue that mistake detection should be **one class** and **online**.



Alessandro Flaborea, Guido D'Amely, Leonardo Plini, Luca Scofano, Edoardo De Matteis, Antonino Furnari, Giovanni Maria Farinella, Fabio Galasso (2024). PREGO: online mistake detection in PRocedural EGOcentric videos . In Conference on Computer Vision and Pattern Recognition (**CVPR 2024**)





Alessandro Flaborea, Guido D'Amely, Leonardo Plini, Luca Scofano, Edoardo De Matteis, Antonino Furnari, Giovanni Maria Farinella, Fabio Galasso (2024). PREGO: online mistake detection in PRocedural EGOcentric videos . In Conference on Computer Vision and Pattern Recognition (**CVPR 2024**)



			Ass	embly101	-0	Epic-tent-O				
	Step Recog.	Step Antic.	Precision	Recall	F1 score	Precision	Recall	F1 score		
One-step memory	Oracle		16.3	30.7	21.3	6.6	26.6	10.6		
BERT [7]	Oracle		78.2	20.0	31.8	75.0	5.6	10.4		
PREGO	Oracle	GPT-3.5	29.2	75.8	42.1	9.9	73.3	17.4		
PREGO	Oracle	LLAMA	30.7	94.0	46.3	10.7	86.7	19.1		
OadTR for MD [35]	<i>OadTR</i> [35]	OadTR [35]	24.3	18.1	20.7	6.7	21.7	10.2		
PREGO	<i>OadTR</i> [35]	LLAMA	22.1	94.2	35.8	9.5	93.3	17.2		
PREGO	MiniRoad [2]	GPT-3.5	16.2	87.5	27.3	4.3	66.6	8.0		
PREGO	MiniRoad [2]	LLAMA	27.8	84.1	41.8	8.6	20.0	12.0		

Alessandro Flaborea, Guido D'Amely, Leonardo Plini, Luca Scofano, Edoardo De Matteis, Antonino Furnari, Giovanni Maria Farinella, Fabio Galasso (2024). PREGO: online mistake detection in PRocedural EGOcentric videos . In Conference on Computer Vision and Pattern Recognition (**CVPR 2024**)





Differentiable Task Graph Learning: Procedural Activity Representation and Online Mistake Detection from Egocentric Videos





Prof. Giovanni Maria Farinella



Antonino Furnari

arXiv

CODE IS AVAILABLE:

https://github.com/fpv-iplab/Differentiable-Task-Graph-Learning

Università Task Graph Learning

We propose the Task Graph Maximum Likelihood (TGML) loss directly supervises the entries of the adjacency matrix Z generating gradients to maximize the probability of edges from past nodes (K3, K1) to the current node (K2), while minimizing the probability of edges from past nodes to future nodes (K4, K5) in a contrastive manner.



CODE IS AVAILABLE: <u>https://github.com/fpv-iplab/Differentiable-Task-Graph-Learning</u>



See paper for the details

CODE IS AVAILABLE:

https://github.com/fpv-iplab/Differentiable-Task-Graph-Learning



Università Online Mistake Detection with Generated Task Graphs

		1	Asser	nbly1	101		EPIC-Tent							
	Avg	Co	Correct			Mistake Avg			Correct			Mistake		
Method	$\overline{F_1}$	F_1	Prec	Rec	F_1	Prec	Rec	$\overline{F_1}$	$\overline{F_1}$	Prec	Rec	F_1	Prec	Rec
Count-Based* 3	26.2	9.5	5.1	86.2	42.9	97.8	27.5	56.6	92.5	92.8	92.2	20.7	20.0	21.4
LLM*	29.3	15.1	8.3	87.2	43.4	96.7	27.9	47.7	86.3	82.4	90.6	9.1	13.3	6.9
MSGI* [<u>35</u>]	33.1	22.7	13.1	84.4	43.5	93.4	28.3	44.5	66.9	51.6	95.2	22.0	73.3	12.9
PREGO* [12]	39.4	32.6	89.7	19.9	46.3	30.7	94.0	32.1	45.0	95.7	29.4	19.1	10.7	86.7
MSG ^{2*} [18]	56.1	63.9	51.5	84.2	48.2	73.6	35.8	54.1	92.9	94.1	91.7	15.4	13.3	18.2
TGT-text (Ours)*	<u>62.8</u>	<u>69.8</u>	56.8	90.6	<u>55.7</u>	84.1	41.7	64.1	93.8	94.1	93.5	34.5	33.3	35.7
DO (Ours)*	75.9	90.2	98.2	83.4	61.6	46.7	90.4	<u>58.3</u>	<u>93.5</u>	94.8	92.4	<u>23.1</u>	20.0	27.3
Improvement*	+19.8	+26.3			+13.4			+7.5	+0.9			+12.5		
Count-Based ⁺ [3]	23.1	2.5	1.3	60.0	43.7	97.8	28.1	38.1	68.3	54.9	90.1	7.9	26.7	4.7
LLM^+	28.1	15.1	7.8	65.5	42.3	89.5	27.7	35.9	61.6	46.7	90.4	10.2	40.0	5.8
MSGI ⁺ [35]	28.4	14.0	7.8	67.9	<u>42.7</u>	90.7	28.0	40.4	59.2	42.9	95.5	21.6	80.0	12.5
PREGO ⁺ [12]	32.5	23.1	68.8	13.9	41.8	27.8	84.1	29.4	41.6	97.9	26.4	17.2	9.5	93.3
MSG ²⁺ [<u>18</u>]	46.2	59.1	51.2	70.0	33.2	44.5	26.5	45.2	67.5	52.4	95.1	22.9	73.3	13.6
TGT-text (Ours) $^+$	<u>53.0</u>	<u>67.8</u>	62.3	74.5	38.2	46.2	32.6	43.8	69.5	55.8	92.1	18.2	53.3	11.0
DO (Ours) ⁺	53.5	78.9	85.0	73.5	28.1	22.5	37.3	46.5	<u>69.3</u>	54.4	95.2	23.7	73.3	14.1
Improvement ⁺	+7.3	+19.8			-5.5			+1.3	+1.2			+1.2		

Table 3: Online mistake detection results. Results obtained with ground truth action sequences are denoted with *, while results obtained on predicted action sequences are denoted with +.



Take home messages

Procedure Understanding opens **many research challenges** as well as the **new and useful applications in different domains to support humans where they live and work** (e.g. Ego-Worker, Ego-Home, ...);





NEXT VISI⁄ິ N

Spin-off of the University of Catania

Synchronization is All You Need: Exocentric-to-Egocentric Transfer for Temporal Action Segmentation with Unlabeled Synchronized Video Pairs







Giovanni Maria Farinella

Poster Session 1 10:30AM TUE 1 OCT

Camillo Quattrocchi

ar

Antonino Furnari Daniele Di Mauro Valerio Giuffrida

C. Quattrocchi, A. Furnari, D. Di Mauro, M. V. Giuffrida, G. M. Farinella (2024). Synchronization is All You Need: Exocentric-to-Egocentric Transfer for Temporal Action Segmentation with Unlabeled Synchronized Video Pairs. European Conference on Computer Vision (ECCV)

Temporal Action Segmentation





Unsupervised Mistake Detection in Egocentric Procedural Video by Detecting Unpredictable Gaze





Michele Mazzamuto

Antonino Furnari Giovanni Maria Farinella

Michele Mazzamuto, Antonino Furnari, and Giovanni Maria Farinella (2024). Eyes Wide Unshut: Unsupervised Mistake Detection in Egocentric Procedural Video by Detecting Unpredictable Gaze. ArXiv: https://arxiv.org/abs/2406.08379

Gaze-Based Mistake Detection



Object-Interaction Anticipation





Despite promising attempts there are <u>still a lot</u> <u>of progress to be done</u> to solve challenges in this research context. Join the effort!



Finanziato dall'Unione europea NextGenerationEU









NEXT VISI ⁽)/N

Spin-off of the University of Catania



Thank you for your attention

Procedure Understanding from Egocentric Videos

Giovanni Maria Farinella

FPV@Image Processing Laboratory - http://iplab.dmi.unict.it/fpv

Next Vision - <u>http://www.nextvisionlab.it/</u>

Department of Mathematics and Computer Science - University of Catania

giovanni.farinella@unict.it - www.dmi.unict.it/farinella