

Towards Visually Intelligent Agents (VIA): a Hybrid Approach

Agnese Chiatti Middle-Late stage PhD Student

@agnese_chiatti achiatti.github.io





Supervisors: Prof. Enrico Motta, Dr. Enrico Daga









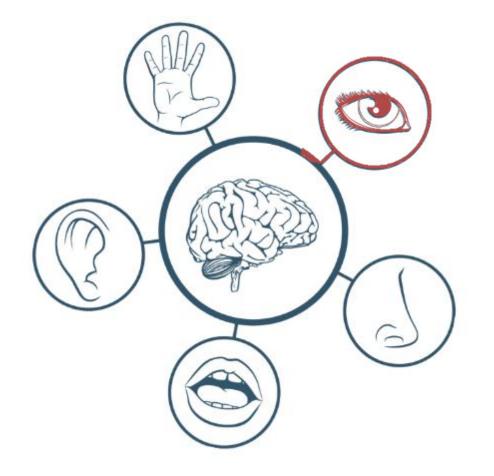
Smart Cities And Robotics

Service Robotics





Research scope





From perception

to **sensemaking**

Visual Intelligence (VI)

"a robot's ability to use its vision system, reasoning components and background knowledge to make sense of the environment." (Chiatti *et al.*, 2020)

Chiatti, A., Motta, E., Daga, E. (2020) *Towards a Framework for Visual Intelligence in Service Robotics: Epistemic Requirements and Gap Analysis*. In Proceedings of KR2020 – Special Session on KR & Robotics.

Background



ML limitations (Marcus, 2018; Pearl, 2018; Parisi et al., 2019)

Promise of hybrid methods (Aditya et al., 2019; Gouidis et al., 2019)

Knowledge-based components at different levels

Post-hoc integration: modularity, isolating the contributing components, transparency but also knowledge reliability and computational overhead to keep in check

Advances in SW and KE: abundance of resources but which to prioritise? (Daruna et al., 2018)

Research Questions



Hypothesis: a hybrid approach (ML+knowledge-based) can improve a robot's performance on tasks that require VI, compared to pure ML.

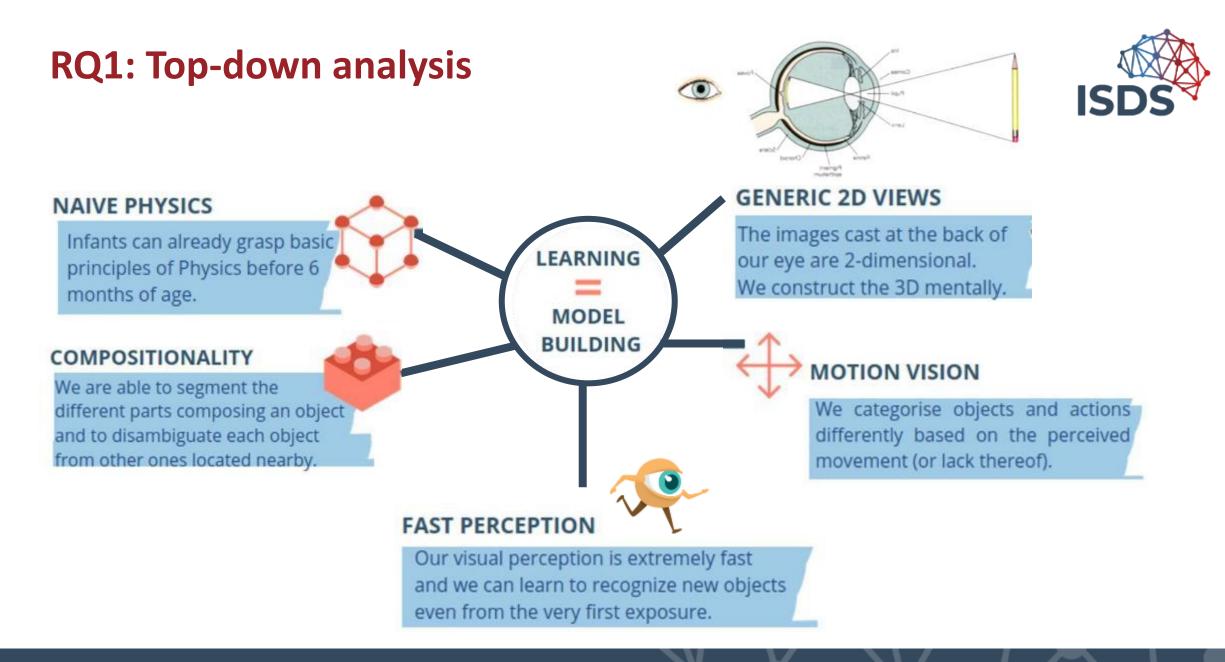
RQ1: What are the epistemic requirements of developing VIA?

RQ2: Which epistemic requirements are **the most important ones**, in our use-case?

RQ3: To what extent do **state-of-the-art KBs** support VIA?

RQ4: To what extent can existing KBs be **repurposed** to support the requirements highlighted in RQ2?

RQ5: To what extent can **a concrete architecture** which integrates the identified types of reasoners be developed?



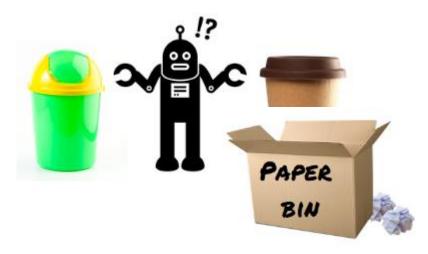
RQ1: Bottom-up analysis



Object recognition pipeline **purely based on ML** (Chiatti et al., 2020)

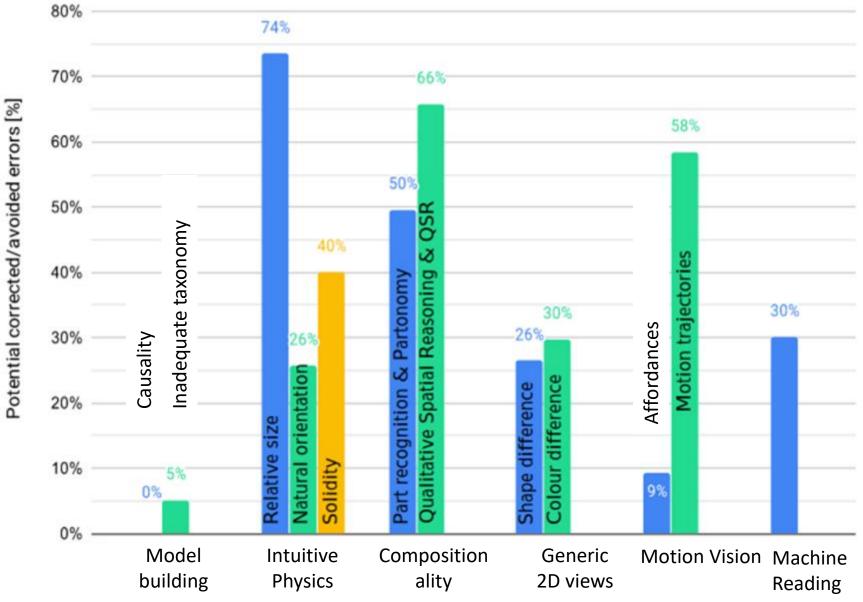
Out of 896 test regions, 272 (31.59%) were mis-classified

Which capability/-ies or knowledge property/-ies can help?



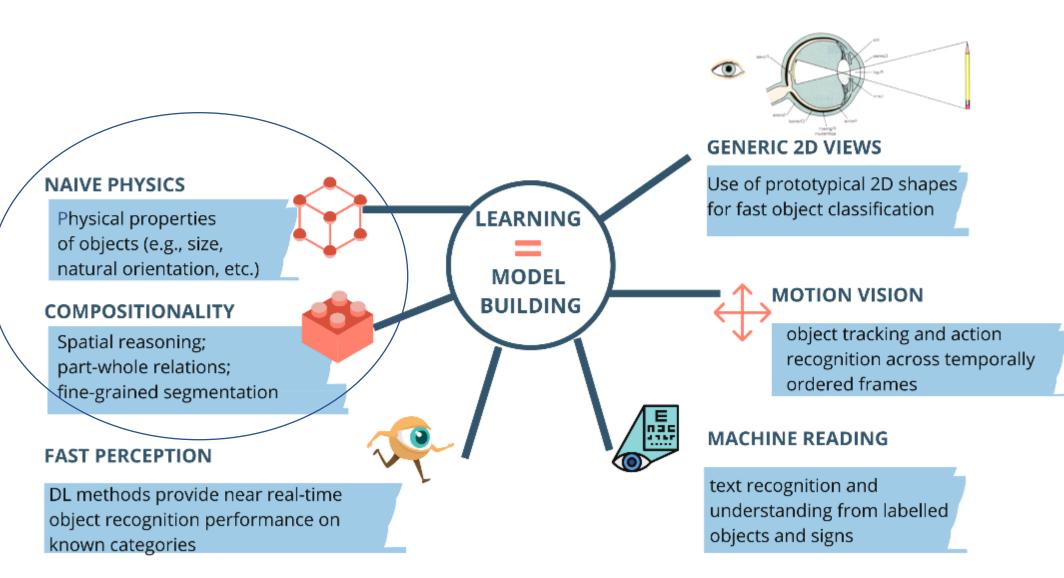
Ground	Predicted	In	tuitive Phys	Spatial	Machine	
truth class	class	Size	Ori.	Solid.	Reason-	Reading
Bin	Mug	T	F	F	Ing T	F
Bin	Mug	T)	F	F	(T)	Т

RQ2: Bottom-up analysis





Visual Intelligence Framework



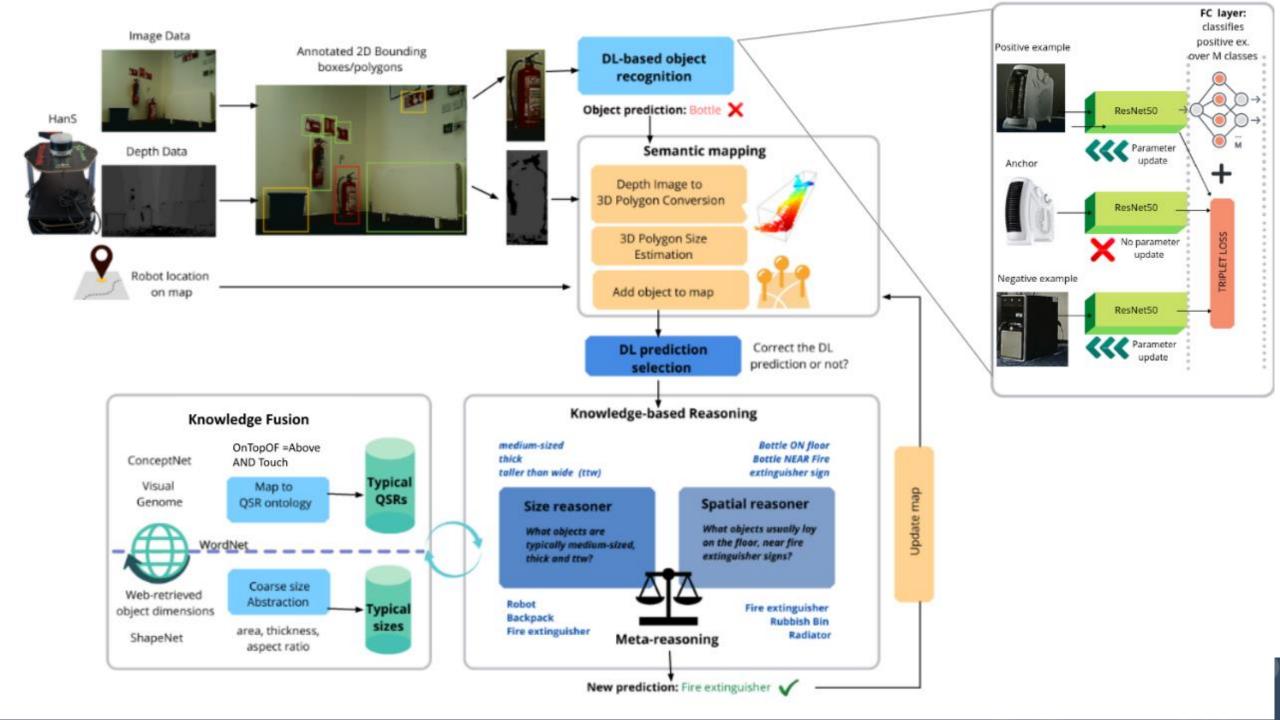


RQ3: KB coverage study



	Knowledge Requirements of Visual Intelligence									
КВ	Hierarchical Taxonomy (linked to WordNet?)	Cause-effect relations	Intuitive Physics Knowledge	Part- whole relations	QSR	Generic 2D views	Object affordances	Motion tra- jectories	Accessibility	
Unified Verb Index (UVI)	yes	0					0		High	
KnowRob/Open-EASE	yes	0	$\bigcirc \bigcirc$	0			\bigcirc	0	Partial	
DBpedia	yes		0	0	0		0		High	
Wikidata	yes		0	0			0		High	
Probase	no	00	0	0					Partial	
NELL	no		0	0	0		00		Adequate	
ConceptNet	yes	00	00	0	0		00		High	
WebChild	yes		0	0	0		00		High	
ATOMIC	no	0			1		0		High	
ASER	no	00			Ì		0		Adequate	
Visual Genome (VG)	yes		00	00	00	0	00		High	
ShapeNet/ PartNet	yes		00	00		00			Partial	

Chiatti, A., Motta, E., Daga, E. (2020) *Towards a Framework for Visual Intelligence in Service Robotics: Epistemic Requirements and Gap Analysis*. In Proceedings of KR2020 – Special Session on KR & Robotics.



Representing sizes

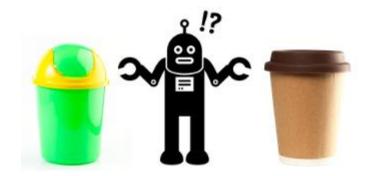
Mid-level representation, robust to contour variance (Long et al., 2016; Zhu et al., 2014; Bagherinezad et al., 2016; Elazar et al., 2019)

Coarse features but also fine-grained enough for categorisation \rightarrow a synoptic representation of size

Object surface area, depth and Aspect Ratio (AR)

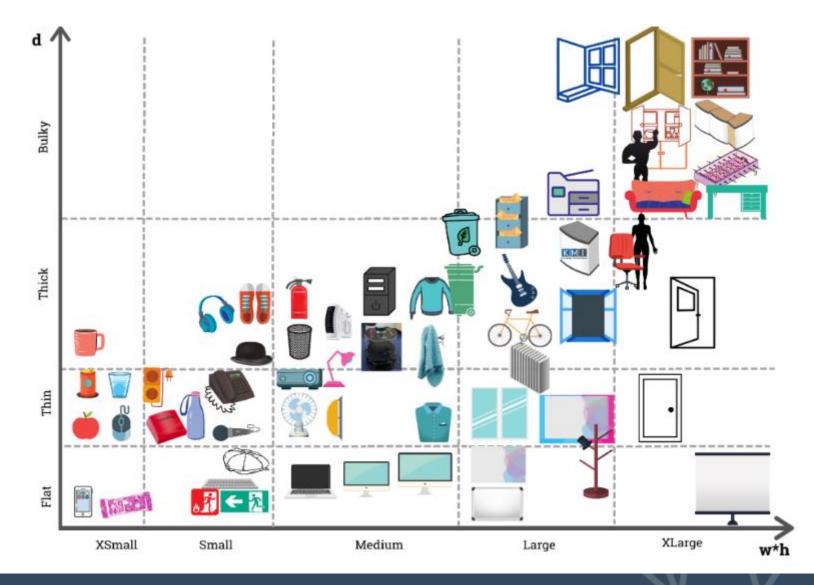
Automatically abstracted from multiple raw size measurements (ShapeNet, Amazon, manually-collected)

Chiatti, A., Motta, E., Daga, E., Bardaro, G. (2021) *Fit to Measure: Reasoning about Sizes for Robust Object Recognition*. In Proceedings of the AAAI-MAKE 2021 Spring Symposium.





Proposed Representation





1¹3

Experimental results: KMi dataset



On a test set of 1414 object regions, annotated w.r.t. 60 reference classes

	Top-1 unweigh. Top-1 weigh.						Top-5 unweigh.		
Method	Р	R	F1	Р	R	F1	P@5 nDCG@5	HR	
N-net [34]	34.0	40.1	31.0	61.5	45.2	47.2	33.1 36.0	63.0	
K-net [34]	39.0	39.9	34.0	68.0	47.9	50.4	$38.5 \ 40.7$	65.1	
Hybrid (area)	39.6	39.5	35.5	65.5	50.3	51.6	41.0 43.1	68.0	
Hybrid (area+flat/non-flat)	41.0	39.3	35.7	65.8	50.1	52.1	$40.5 \ 42.8$	65.8	
Hybrid (area+thickness)	44.5	38.9	38.6	65.0	51.4	53.9	$41.8 \ 44.1$	68.5	
Hybrid (area+flat/non-flat+AR)	42.9	38.8	36.6	68.9	49.1	52.9	$39.9 \ 42.0$	66.3	
Hybrid (area+thickness+AR)	47.2	39.1	40.0	69.1	51.4	55.4	$41.6 \ 43.9$	68.4	

Experimental results: Amazon 2017 Image Matching dataset



On a test set of 562 images, annotated w.r.t. 61 classes (41 known + 20 novel)

	Top-1 accuracy			Top-5 unweighted		
Method	Known	Novel	Mixed	P@5	nDCG@5	HR
N-net [34]	56.8	82.1	64.6	61.9	62.7	72.6
K-net [34]	99.7	29.5	78.1	73.7	75.0	82.4
Hybrid (area)	94.7	71.7	87.6	82.6	84.1	89.7
Hybrid (area $+ $ flat/non-flat)	94.5	71.7	87.5	82.5	84.0	89.7
Hybrid (area + thickness)	81.7	39.3	68.7	64.6	65.8	70.1

Representing Qualitative Spatial Relations (QSR)

Several different AI formalisms (Cohn& Renz., 2008)

Semantic mapping methods & GIS technologies (Kostavelis & Gasteratos, 2015)

Best of both worlds: defining a mapping (Deeken et al., 2018)

But also account for everyday language use to describe spatial relations (Landau & Jackendoff, 1993)

Which is also reflected in KBs like Visual Genome, ConceptNet, SpatialSense, etc.

Broader impact on HRI (Sarthou et al., 2019; Sisbot & Connell, 2019)





Representing Qualitative Spatial Relations (QSR)



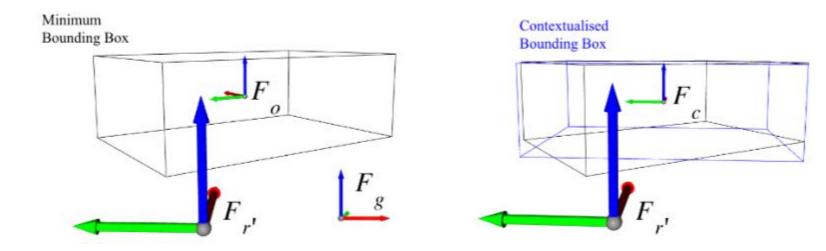


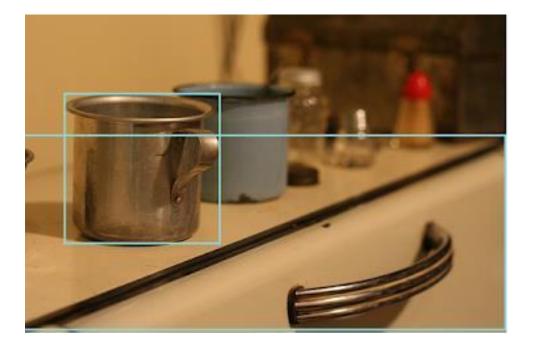
Figure and reference but also dependent on the observer

Chiatti, A., Bardaro, G., Motta, E., Daga, E. (2021) *Commonsense Spatial Reasoning for Visually Intelligent Agents*. <u>https://arxiv.org/abs/2104.00387</u>

Extracting spatial statistics from VG



Different spatial uses of "on" predicate: on top of, leans on, affixed on





Img source: https://visualgenome.org/

Conclusion and next steps



Before we can delegate tasks to robots, we need to enhance their sensemaking capabilities.

Contributed a framework of requirements for VIA and guidelines on priorities

Preliminary evidence that **knowledge** of the typical size and spatial relations between objects **integrated in post-processing can significantly augment ML**.

Evaluation of the spatial reasoning component in progress

Meta-reasoning: study of cases of agreement/disagreement between reasoners

Completing assessment of utility to support **decision-making scenarios**



Thank you! Q&A



robots.kmi.open.ac.uk



https://github.com/kmi-robots



@agnese_chiatti @isdsou







Linked with common sense

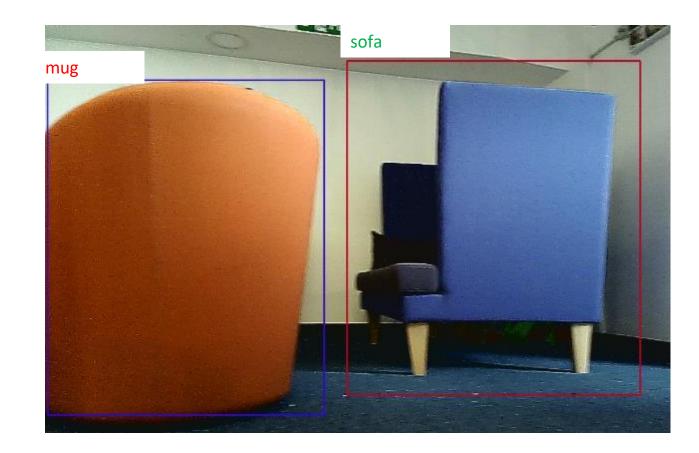
Experiential

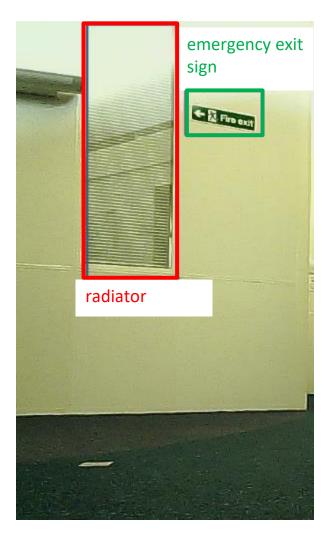
Cross-sectional (Space, Time, Physics,...)



ML limitations







22