

# Project Part 1

## Computational Visual Perception (CompVP)

Bernhard Egger, Andreas Kist, Patrick Krauß, Tim Weyrich

- 7,5 ECTS
  - 5+2,5 ECTS
  - You need both!
  - You can't get only 5 or only 2,5 ECTS

- How much of vision do SotA generative video models “solve”?
  - How well do they and other models work for corner cases of vision?
-

## Video models are zero-shot learners and reasoners

Thaddäus Wiedemer<sup>\*1</sup>, Yuxuan Li<sup>1</sup>, Paul Vicol<sup>1</sup>, Shixiang Shane Gu<sup>1</sup>, Nick Matarese<sup>1</sup>, Kevin Swersky<sup>1</sup>,  
Been Kim<sup>1</sup>, Priyank Jaini<sup>\*1</sup> and Robert Geirhos<sup>\*1</sup>

<sup>1</sup>Google DeepMind

The remarkable zero-shot capabilities of Large Language Models (LLMs) have propelled natural language processing from task-specific models to unified, generalist foundation models. This transformation emerged from simple primitives: large, generative models trained on web-scale data. Curiously, the same primitives apply to today's generative video models. Could video models be on a trajectory towards general-purpose *vision* understanding, much like LLMs developed general-purpose *language* understanding? We demonstrate that Veo 3 can solve a broad variety of tasks it wasn't explicitly trained for: segmenting objects, detecting edges, editing images, understanding physical properties, recognizing object affordances, simulating tool use, and more. These abilities to perceive, model, and manipulate the visual world enable early forms of visual reasoning like maze and symmetry solving. Veo's emergent zero-shot capabilities indicate that video models are on a path to becoming unified, generalist vision foundation models.

Project page: <https://video-zero-shot.github.io/>

### 1. Introduction

29 Sep 2025



## Illusory Motion Reproduced by Deep Neural Networks Trained for Prediction

Eiji Watanabe<sup>1,2\*</sup>, Akiyoshi Kitaoka<sup>3</sup>, Kiwako Sakamoto<sup>4,5</sup>, Masaki Yasugi<sup>1</sup> and Kenta Tanaka<sup>6</sup>

<sup>1</sup> Laboratory of Neurophysiology, National Institute for Basic Biology, Okazaki, Japan, <sup>2</sup> Department of Basic Biology, The Graduate University for Advanced Studies (SOKENDAI), Miura, Japan, <sup>3</sup> Department of Psychology, Ritsumeikan University, Kyoto, Japan, <sup>4</sup> Department of Physiological Sciences, The Graduate University for Advanced Studies (SOKENDAI), Miura, Japan, <sup>5</sup> Division of Integrative Physiology, National Institute for Physiological Sciences (NIPS), Okazaki, Japan, <sup>6</sup> Sakura Research Office, Wako, Japan

The cerebral cortex predicts visual motion to adapt human behavior to surrounding objects moving in real time. Although the underlying mechanisms are still unknown, predictive coding is one of the leading theories. Predictive coding assumes that the brain's internal models (which are acquired through learning) predict the visual world at all times and that errors between the prediction and the actual sensory input further refine

- Project Part 1: Play around with video diffusion models and reproduce paper (close to Paper 1)
  - Project Part 2: Experiment further with illusions, optical flow and depth estimation (close to paper 2)
  - Project Part 3: Further Evaluation based on results of Part 1 and Part 2 (close to paper 1 and 2)
-

- Submission via Studon course:  
“Computational Visual Perception”
  - Submission has to be in the exact format
  - Three strict project deadlines
    - November 21st,
    - January 2nd (feel free to submit early),
    - February 5th
-

- Pass/Fail for each of the 3 parts,
- You need to pass 2 out of 3 parts
- The best solutions for each part will be released ~ 1 week after the deadline, to enable others to continue with the best solution of another team



# How to pass

---

- Scope of project ~ 150 hours per student
  - Teams of 1-3 students
  - Steps can be performed in new group
  - If you are looking for a group, please stay after the class and talk to people who also stay
  - If you are looking for a group and can only join virtually, please use the forum in “Computational Visual Perception” to team up
  - Finding a group is your responsibility
-

# Part 1 Task 1

- Choose 3 cases you would like to reproduce
- Reproduce those cases with the model of your choice
  - Model examples:
    - Veo 3 (closed)
    - Wan 2.2 (open)
    - <https://lmarena.ai/leaderboard/text-to-video>

**Text-to-Video Arena**  
Compare models according to their ability to generate videos based on the given prompt.  
Generate videos and vote in the [Discord server](#).

Last Updated: Oct 19, 2025    Total Votes: 88,217    Total Models: 28

Overall    Search by model name...

Rank (UB)	Model	Score	95% CI (s)	Votes	Organization	License
1	veo-3.1-audio	1404	±20	1,305	Google	Proprietary
1	veo-3.1-fast-audio	1395	±19	1,334	Google	Proprietary
1	sora-2-pro	1365	±21	2,459	OpenAI	Proprietary
2	veo-3-fast-audio	1368	±13	21,532	Google	Proprietary
3	veo-3-audio	1354	±13	15,107	Google	Proprietary
6	sora-2	1318	±17	2,697	OpenAI	Proprietary
7	veo-3-fast	1261	±13	11,304	Google	Proprietary
7	veo-3	1247	±14	11,023	Google	Proprietary
8	klimg-2.5-turbo-1080p	1223	±17	1,627	KlingAI	Proprietary

# Part 1 Task 2

---

- Make your own variant for each of the 3 cases you did choose
- Take your own image with a camera
  - No AI generated image
  - No images from the internet



# Part 1 Task 2 voluntary part

---

- Collect ground truth (if applicable)



# Part 1 Task 3

- Come up with one own task
- If you need inspiration:

## A Definition of AGI








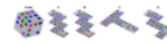




Dan Hendrycks<sup>1</sup>, Dawn Song<sup>2</sup>, Christian Szegedy<sup>3</sup>, Honglak Lee<sup>4,5</sup>, Yarin Gal<sup>6</sup>, Erik Brynjolfsson<sup>7</sup>, Sharon Li<sup>8</sup>, Andy Zou<sup>1,9,10</sup>, Lionel Levine<sup>11</sup>, Bo Han<sup>12</sup>, Jie Fu<sup>13</sup>, Ziwei Liu<sup>14</sup>, Jinwoo Shin<sup>15</sup>, Kimin Lee<sup>15</sup>, Mantas Mazeika<sup>1</sup>, Long Phan<sup>1</sup>, George Ingebrechtsen<sup>1</sup>, Adam Khoja<sup>1</sup>, Cihang Xie<sup>16</sup>, Olawale Salaudeen<sup>17</sup>, Matthias Hein<sup>18</sup>, Kevin Zhao<sup>19</sup>, Alexander Pan<sup>2</sup>, David Duvenaud<sup>20,21</sup>, Bo Li<sup>22</sup>, Steve Omohundro<sup>23</sup>, Gabriel Alfour<sup>24</sup>, Max Tegmark<sup>17</sup>, Kevin McGrew<sup>25</sup>, Gary Marcus<sup>26</sup>, Jaan Tallinn<sup>27</sup>, Eric Schmidt<sup>17</sup>, Yoshua Bengio<sup>28,29</sup>

<sup>1</sup>Center for AI Safety <sup>2</sup>University of California, Berkeley <sup>3</sup>Morph Labs  
<sup>4</sup>University of Michigan <sup>5</sup>LG AI Research <sup>6</sup>University of Oxford <sup>7</sup>Stanford University  
<sup>8</sup>University of Wisconsin–Madison <sup>9</sup>Gray Swan AI <sup>10</sup>Carnegie Mellon University  
<sup>11</sup>Cornell University <sup>12</sup>Hong Kong Baptist University <sup>13</sup>HKUST  
<sup>14</sup>Nanyang Technological University <sup>15</sup>KAIST <sup>16</sup>University of California, Santa Cruz  
<sup>17</sup>Massachusetts Institute of Technology <sup>18</sup>University of Tübingen <sup>19</sup>University of Washington  
<sup>20</sup>University of Toronto <sup>21</sup>Vector Institute <sup>22</sup>University of Chicago  
<sup>23</sup>Beneficial AI Research <sup>24</sup>Conjecture <sup>25</sup>Institute for Applied Psychometrics  
<sup>26</sup>New York University <sup>27</sup>CSER <sup>28</sup>Université de Montréal <sup>29</sup>LawZero

## 10 Visual Processing (V)

### Visual Processing (V)

The ability to analyze and generate natural and unnatural images and videos

Perception	Visual Generation	Visual Reasoning	Spatial Scanning
<p>The ability to process and interpret visual inputs from images and videos</p> <p><b>Image Recognition</b></p>  <ul style="list-style-type: none"> <li>• "What does this image depict?"</li> </ul> <p><b>Image Captioning</b></p>  <ul style="list-style-type: none"> <li>• "Create descriptive caption for this image."</li> </ul> <p><b>Image Anomaly Detection</b></p>  <ul style="list-style-type: none"> <li>• "Which is the odd one out?"</li> </ul> <p><b>Clip Captioning</b></p>  <ul style="list-style-type: none"> <li>• "What happens in this video?"</li> </ul> <p><b>Video Anomaly Detection</b></p>  <ul style="list-style-type: none"> <li>• "Is this physically plausible?"</li> </ul>	<p>The ability to synthesize images and short videos</p> <p><b>Simple Natural Images</b></p> <ul style="list-style-type: none"> <li>• "Generate an image of a golden retriever playing in a park."</li> </ul> <p><b>Complicated Images</b></p> <ul style="list-style-type: none"> <li>• "Generate a diagram showing the process of photosynthesis."</li> </ul> <p><b>Simple Natural Videos</b></p> <ul style="list-style-type: none"> <li>• "Generate a short video of somebody typing on a keyboard."</li> </ul>	<p>The ability to understand and inferences about the images</p> <p><b>Gestalt</b></p>  <ul style="list-style-type: none"> <li>• "Identify the picture."</li> </ul> <p><b>Mental Rotation</b></p>  <ul style="list-style-type: none"> <li>• "Which shape on the right is the same as the shape on the left?"</li> </ul> <p><b>Mental Folding</b></p>  <ul style="list-style-type: none"> <li>• "Which net, when folded, cannot form the cube?"</li> </ul> <p><b>Embodied Reasoning</b></p>  <ul style="list-style-type: none"> <li>• "Which trajectories should the zipper follow to zip the suitcase?"</li> </ul> <p><b>Chart and Figure Reasoning</b></p>  <ul style="list-style-type: none"> <li>• "What is the lowest labeled tick on the y-axis?"</li> </ul>	<p>The ability to understand and inferences about the images</p> <ul style="list-style-type: none"> <li>• "Find the path to the center of this maze."</li> </ul>  <ul style="list-style-type: none"> <li>• "Count the people in the picture."</li> </ul> 

**Assessment Details.** See Appendix H for further details on how to assess visual processing capabilities concretely.

- Or from the lecture

- Retry everything with a different model
  - Older/newer version
  - Open Source vs. commercial
  - ...

- In one of the models you tried: deactivate prompt rewriting
  - Wan 2.2. definitely possible
  - Veo over API
- (that needs to be taken into account when choosing models)

- Summarize results in form of a presentation
- (instead of a report, there will be no actual presentation)



# Part 1 deliverables

- Per project team 1 single pptx file (or powerpoint compatible)
- The pptx file contains:
  - Title slide with all team members names
  - 3x1 slide for each reproduced result
    - Titel: used model
    - Left original video, right reproduced video
    - Bottom prompts
  - 3x1 slide for your own variant of reproduced task
    - Titel: used model
    - Left original input image and your input image
    - Right generated videos
    - Bottom prompts
  - Generate one task yourself, same format
  - 2 x 3x1 +1 slide: All results above with a different model, same format
  - 2 x 3x1 +1 slide: All results above also without prompt rewriting, same format
  - 1 slide summarizing what you observed (text)
  - Failure cases (if any)
- All videos must be embedded in the file (no external links)
- (ideally all videos are set to play automatically)

- In Studon course :  
“Computational Visual Perception”
- You will upload up to ~1GB (don't!)  
Plan in internet speed, upload at university

If you run into issues uploading, you send an md5 hash of your zip file **before** the deadline and you provide an alternative download link within 24h

- Correct format
  - 3 results reproduced
  - 3 own variants videos
  - Own task
  - All 7 videos above with a different model
  - All above 7 videos without prompt rewriting
  - Summary text
- 
- Selected solution: all points fulfilled to full satisfaction
  - Pass: at most one bullet point from the above missing
- 
- Plagiarism will have serious consequences
-

- Project Part 1: Play around with video diffusion models and reproduce paper (close to Paper 1)
  - Project Part 2: Experiment further with illusions, optical flow and depth estimation (close to paper 2)
  - Project Part 3: Further Evaluation based on results of Part 1 and Part 2 (close to paper 1 and 2)
-

- You can ask questions in the forum  
“Computational Visual Perception”
  - You come with concrete questions
  - I’ll open a thread in the forum, where you can respond till Thursday each week if you want to meet
  - I’ll distribute time slots each Friday
  - No guarantee for any responses on the day of the deadline
-

- 7,5 ECTS
  - 5+2,5 ECTS
  - You need both!
  - You can't get only 5 or only 2,5 ECTS

# Don't start late

## PROCRASTINATION



jp