



Università
di Catania

CVPR
JUNE 3-7, 2026



DENVER
COLORADO



UNIVERSITY OF
BATH

ViterbiPlanNet: Injecting Procedural Knowledge via Differentiable Viterbi for Planning in Instructional Videos

Highlight



Luigi Seminara

Department of Mathematics and Computer Science

University of Catania

Last Year PhD Student



Davide Moltisanti
University of Bath



Antonino Furnari
University of Catania



@luseminara.bsky.social



@Gigii_Gii



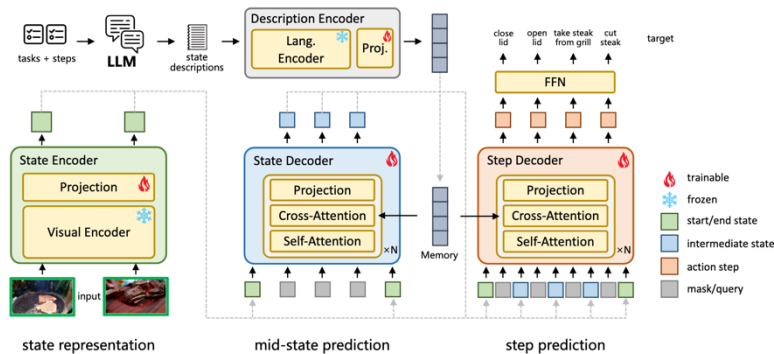
luigi-seminara



Prior Works

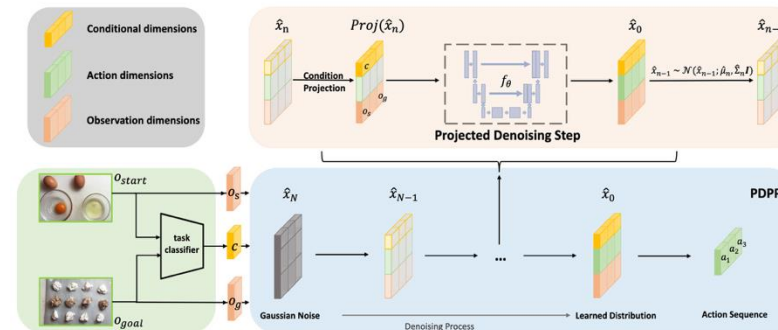
SCHEMA: STATE CHANGES MATTER FOR PROCEDURE PLANNING IN INSTRUCTIONAL VIDEOS

Yulei Niu¹ Wenliang Guo¹ Long Chen² Xudong Lin¹ Shih-Fu Chang¹
¹Columbia University ²The Hong Kong University of Science and Technology
 yn.yuleiniu@gmail.com



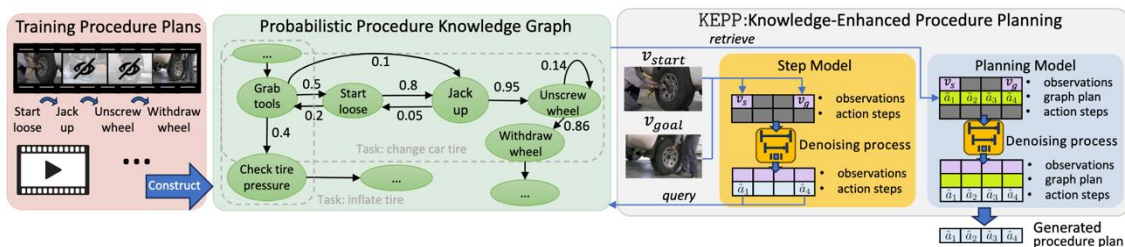
PDPP: Projected Diffusion for Procedure Planning in Instructional Videos

Hanlin Wang¹ Yilu Wu¹ Sheng Guo³ Limin Wang^{1, 2, *}
¹State Key Laboratory for Novel Software Technology, Nanjing University, China
²Shanghai AI Lab, China ³MYbank, Ant Group, China



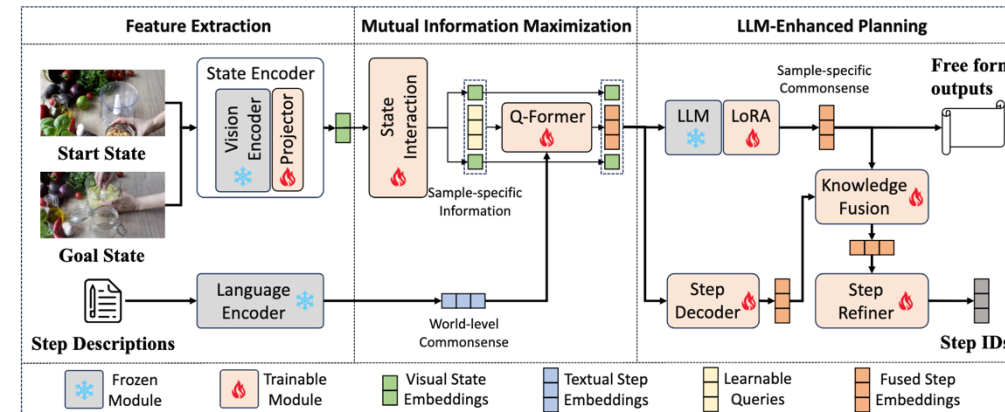
Why Not Use Your Textbook? Knowledge-Enhanced Procedure Planning of Instructional Videos

Kumaranage Ravindu Yasanghe¹ Honglu Zhou² Malitha Gunawardhana^{1,3}
 Martin Renqiang Min² Daniel Harari⁴ Muhammad Haris Khan¹
¹Mohamed bin Zayed University of Artificial Intelligence, ²NEC Laboratories, USA,
³University of Auckland, ⁴Weizmann Institute of Science



PlanLLM: Video Procedure Planning with Refinable Large Language Models

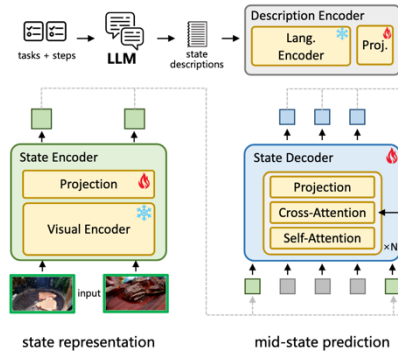
Dejie Yang¹, Zijing Zhao¹, YangLiu^{1,2,*}
¹Wangxuan Institute of Computer Technology, Peking University
²State Key Laboratory of General Artificial Intelligence, Peking University
 {vdi.ziizhzhao}@stu.pku.edu.cn, yangliu@pku.edu.cn



Prior Works

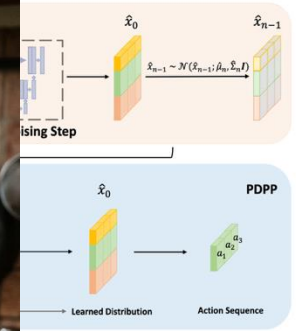
SCHEMA: STATE CHANGES MATTER FOR PROCEDURE PLANNING IN INSTRUCTIONAL VIDEOS

Yulei Niu¹ Wenliang Guo¹ Long Chen² Xudong Lin¹ Shih-Fu Chang¹
¹Columbia University ²The Hong Kong University of Science and Technology
 yn.yuleiniu@gmail.com



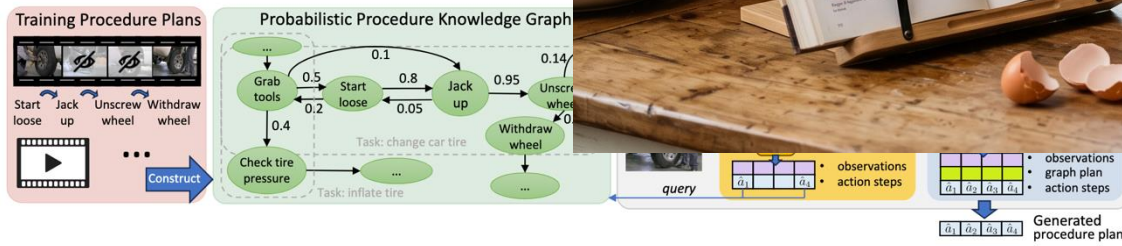
PDPP: Projected Diffusion for Procedure Planning in Instructional Videos

Hanlin Wang¹ Yilu Wu¹ Sheng Guo³ Limin Wang^{1, 2, ✉}
¹State Key Laboratory for Novel Software Technology, Nanjing University, China
²Ant Group, China



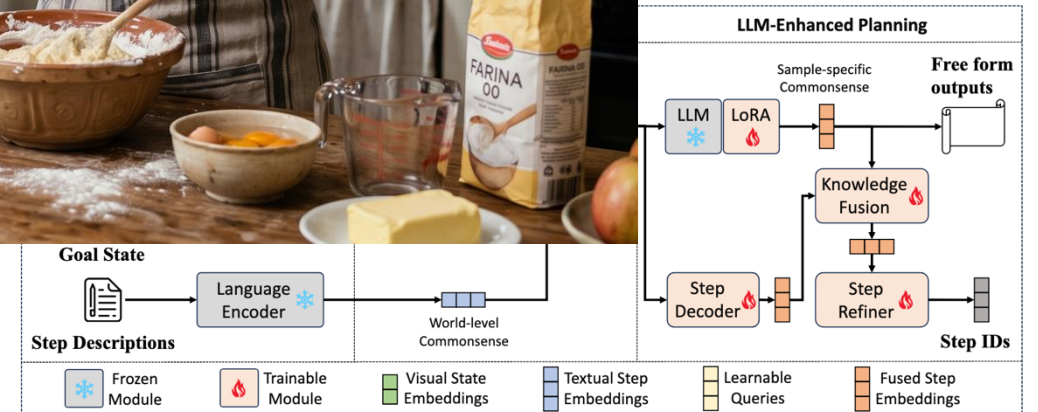
Why Not Use Your Textbook? Knowledge Injection for Instructional Video Planning

Kumarange Ravindu Yasanghe¹ Martin Renqiang Min² Daniel Hoi³
¹Mohamed bin Zayed University of Artificial Intelligence
²University of Auckland, ³University of Toronto

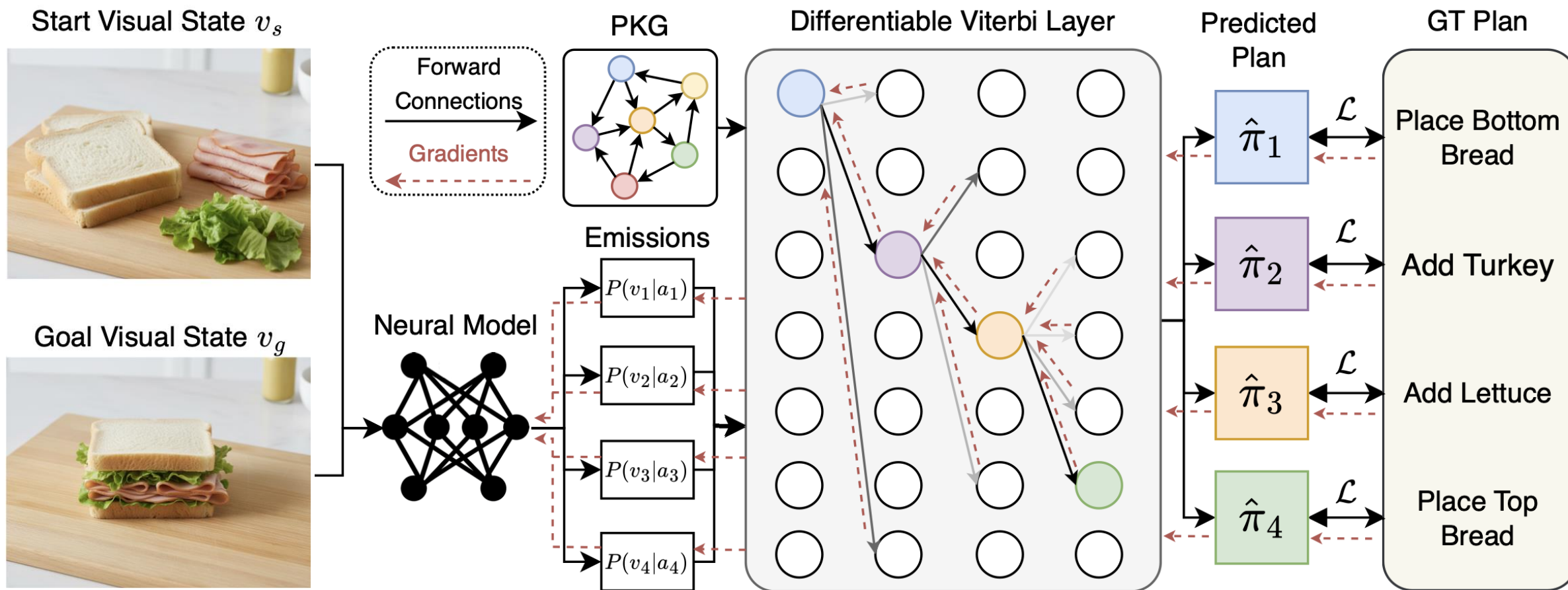


Diffusion-based Large Language Models for Instructional Video Planning

Yi Liu^{1, 2, *} Zhenyuan Liu¹
¹Peking University
²State Key Laboratory of General Artificial Intelligence, Peking University
 liuyi@pku.edu.cn

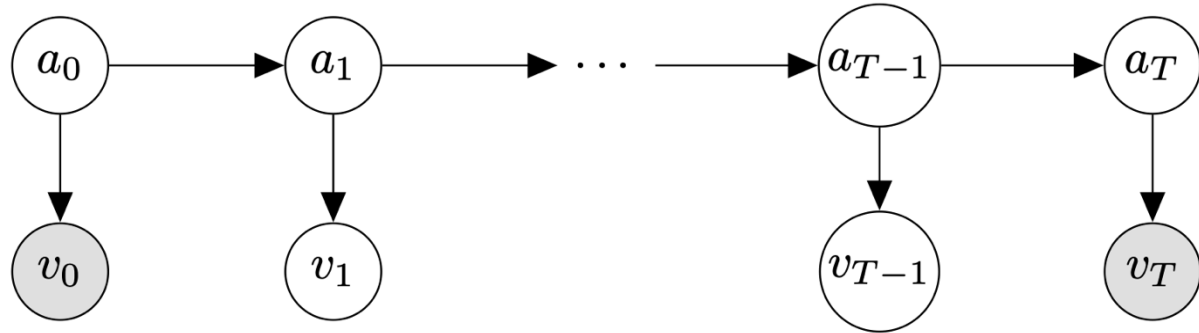


Using Procedural Graphs



Cognitive Offloading to the Graph Available both at Training and Inference

Probabilistic Framework



$$P(a_{0:T}, v_{0:T}) = P(a_0)P(v_0|a_0) \prod_{t=1}^T P(a_t|a_{t-1})P(v_t|a_t)$$

$$\pi^* = \arg \max_{\pi = a_{1:T} \in \mathcal{K}^T} \prod_{t=1}^T \underbrace{P(a_t|a_{t-1})}_{\text{Transition}} \underbrace{P(v_t|a_t)}_{\text{Emission}}$$

Viterbi Algorithm

$$\text{S-max}(\mathbf{x}) = \log \left(\sum_{k=1}^N \exp(x_k - m) \right) + m$$

$$\text{S-argmax}(\mathbf{x})_k = \frac{\exp(x_k - m)}{\sum_{j=1}^N \exp(x_j - m)}$$

Differentiable Relaxations

Algorithm 1 Differentiable Viterbi Layer (DVL)

Input: Emission probabilities $b \in [0, 1]^{T \times N}$, transition probabilities $\omega(i, j)$ from PKG

Output: Soft plan distribution $\tilde{\pi} \in [0, 1]^{T \times N}$

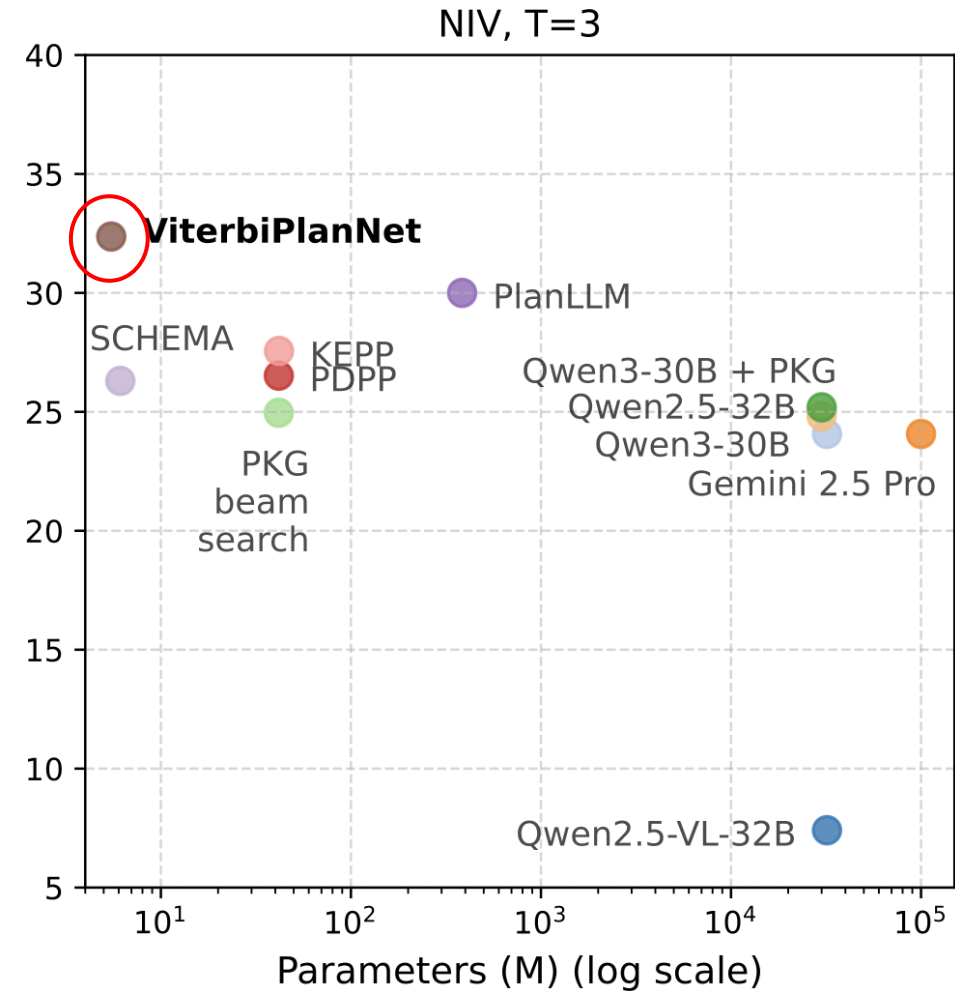
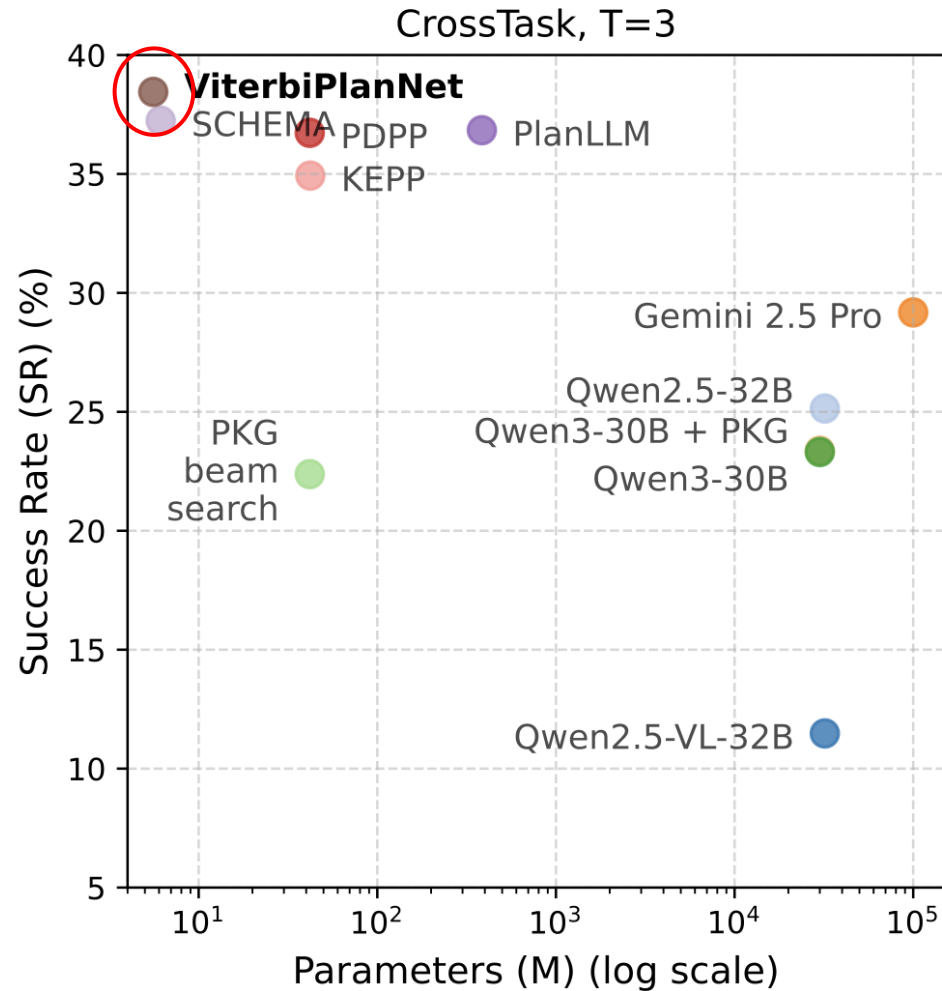
```

1: // 1. Initialization (t=1)
2: for j = 1 to N do
3:    $\delta_1(j) \leftarrow b[1, j]$  ▷ Initial state scores
4: end for
5: // 2. Forward Pass: Differentiable Recursion (t > 1)
6: for t = 2 to T do
7:   for j = 1 to N do
8:     // Compute predecessor scores
9:     for i = 1 to N do
10:       $s_{i \rightarrow j}^{(t)} \leftarrow \delta_{t-1}(i) \cdot \omega(i, j)$ 
11:    end for
12:    // Update state scores using smooth max
13:     $\delta_t(j) \leftarrow b[t, j] \cdot \text{S-max}(\{s_{i \rightarrow j}^{(t)}\}_{i=1}^N)$ 
14:    // Compute soft backpointer distribution
15:     $\psi_t(j, \cdot) \leftarrow \text{S-argmax}(\{s_{i \rightarrow j}^{(t)}\}_{i=1}^N)$ 
16:  end for
17: end for
18: // 3. Backward Pass: Recursive Composition
19:  $\tilde{\pi}_T \leftarrow \text{S-argmax}(\delta_T)$  ▷ Distribution at the final horizon
20: for t = T - 1 down to 1 do
21:    $\tilde{\pi}_t \leftarrow \sum_{j=1}^N \tilde{\pi}_{t+1}(j) \cdot \psi_{t+1}(j, \cdot)$  ▷ Recursive backtrace
22: end for
23: return  $\tilde{\pi}$ 

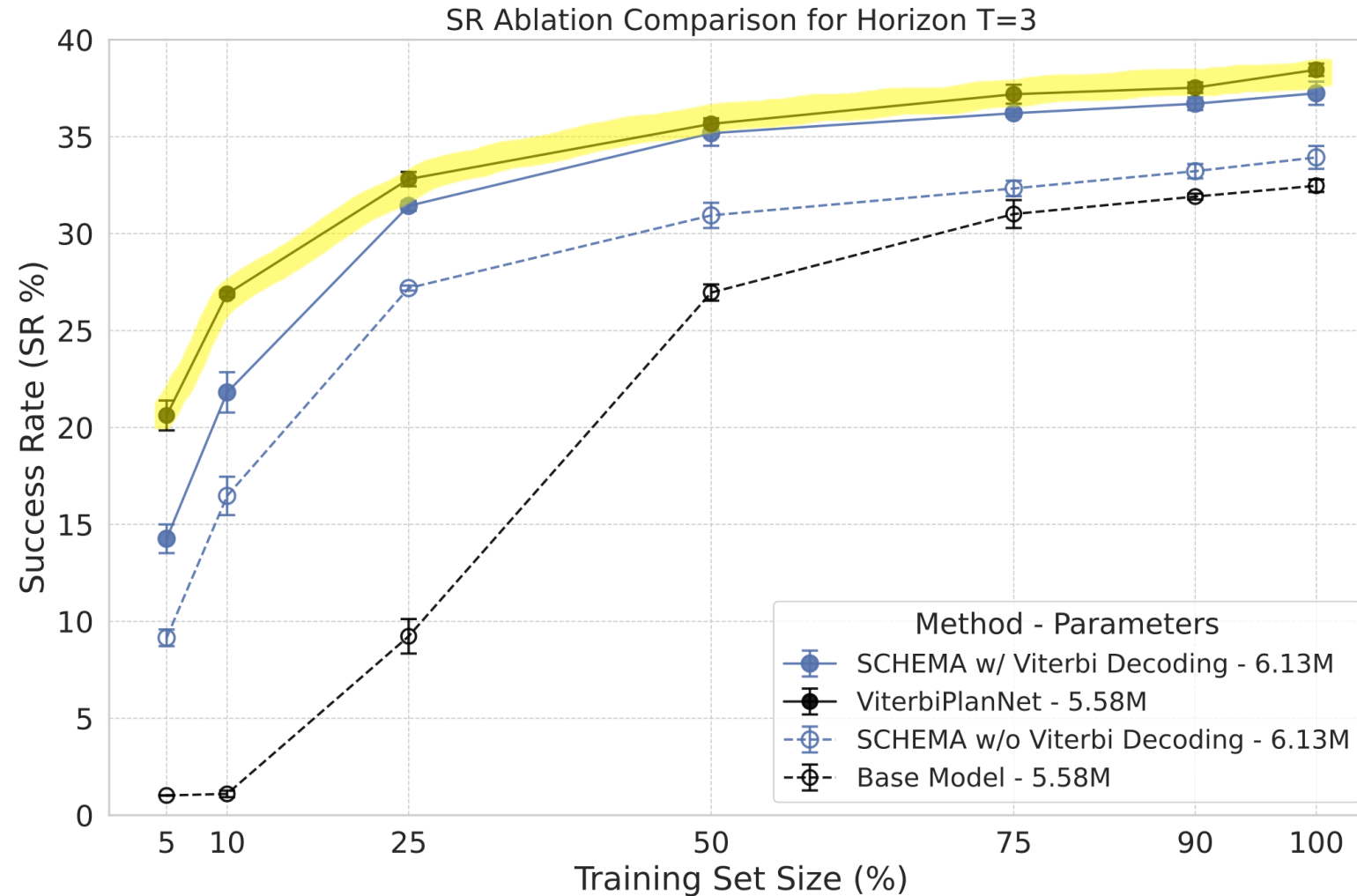
```

Allows Gradients to Pass

Better Performance with Fewer Parameters

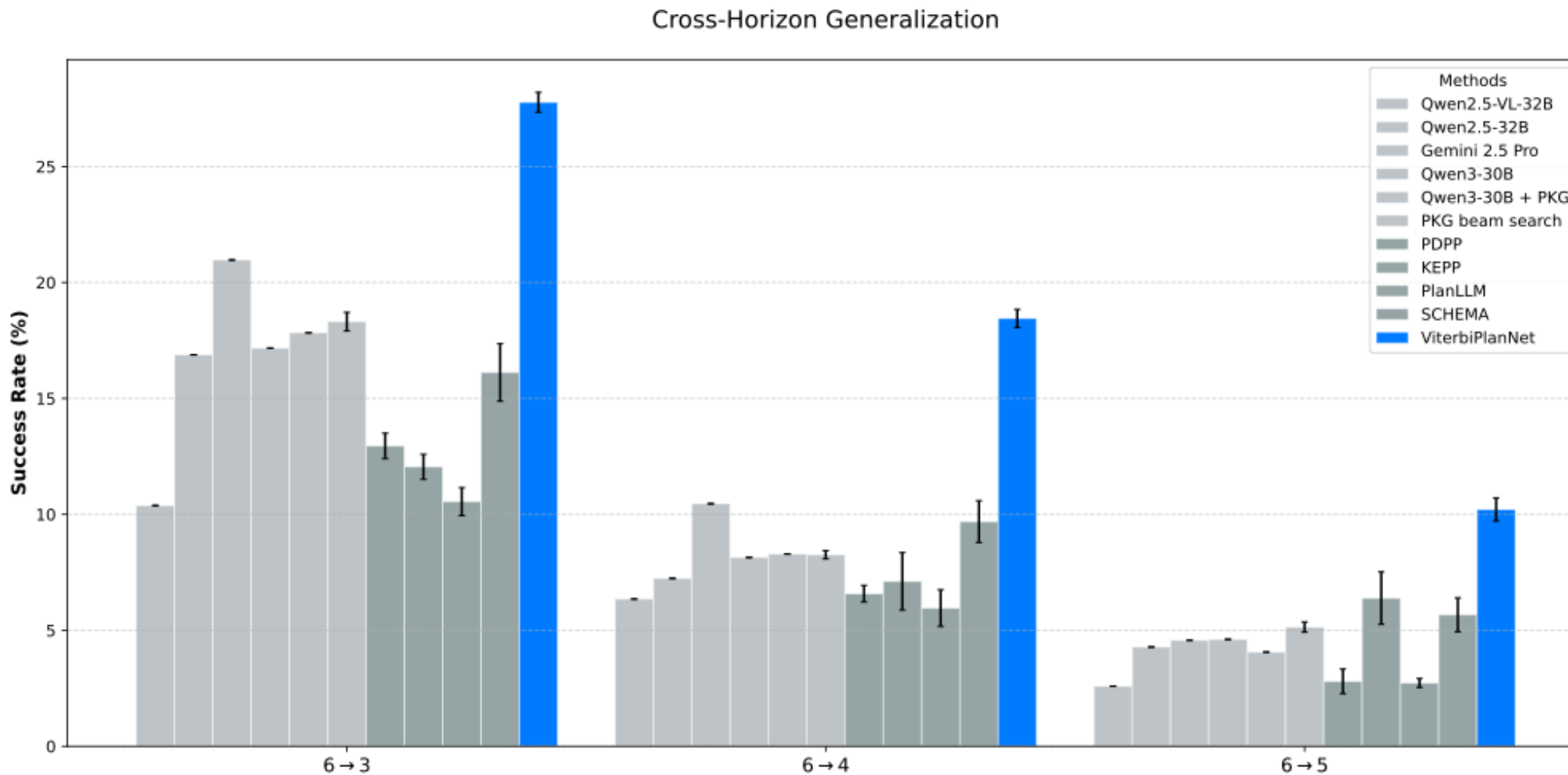


Sample Efficiency



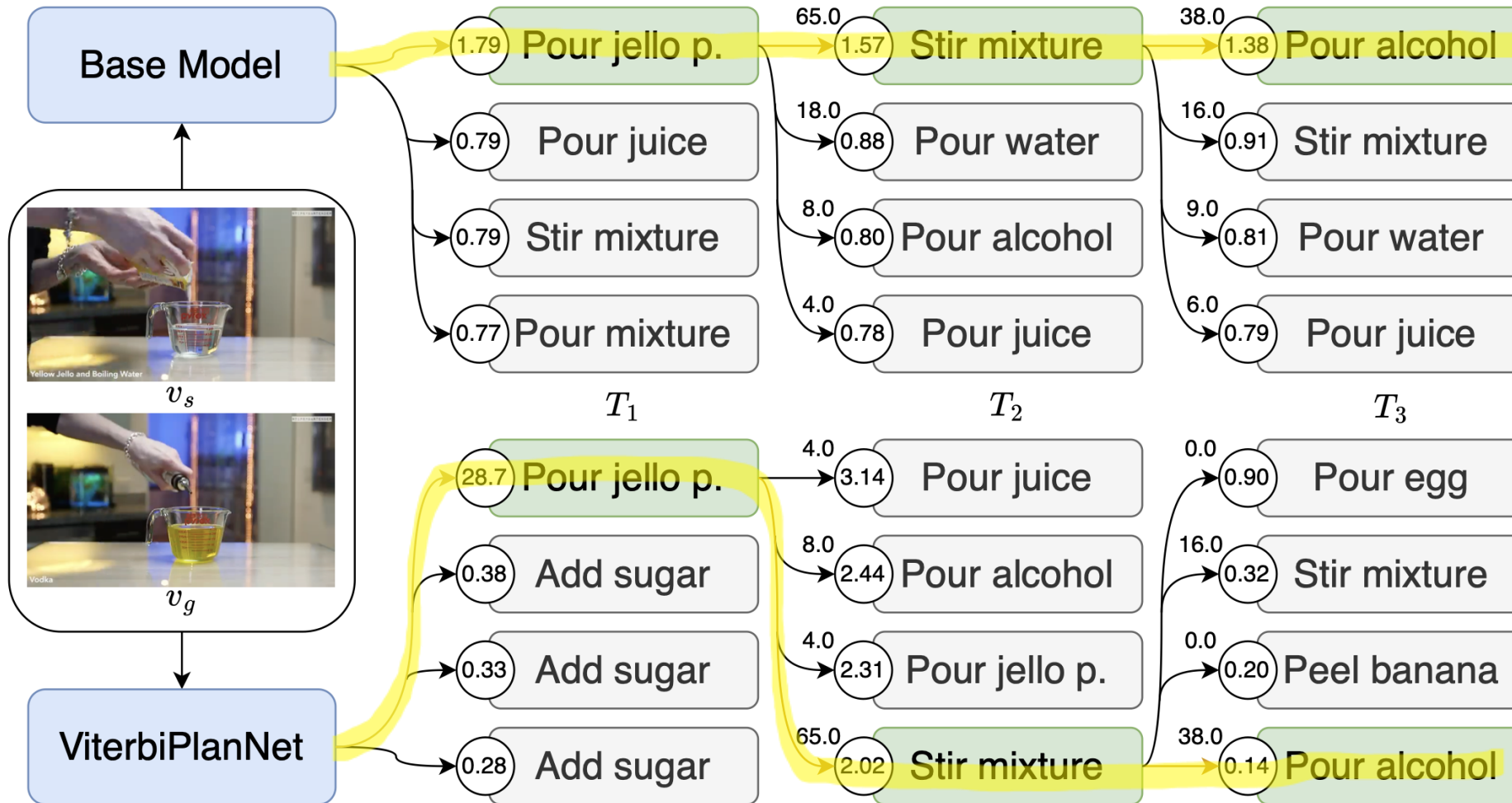
Good performance with fewer training examples.

Cross-Horizon Generalization



Generalizes when trained on a longer horizon and tested on a shorter one.

Qualitative Example

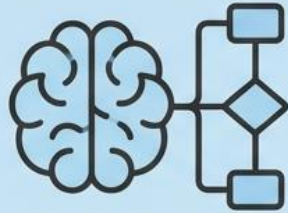


The base model learns to memorize the PKG.

ViterbiPlanNet just has to get a reasonable sorting, relying on the graph for correction.

Easier Cognitive Task

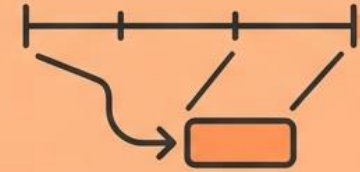
Summary



Avoid memorizing procedural knowledge.



Good performance with fewer training examples.



Generalizes from long to short horizons.



Università
di Catania

CVPR
JUNE 3-7, 2026



DENVER
COLORADO



UNIVERSITY OF
BATH

ViterbiPlanNet: Injecting Procedural Knowledge via Differentiable Viterbi for Planning in Instructional Videos



Highlight

Luigi Seminara

Department of Mathematics and Computer Science

University of Catania

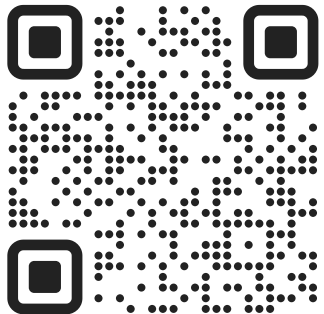
Last Year PhD Student



Davide Moltisanti
University of Bath



Antonino Furnari
University of Catania



@luseminara.bsky.social



@Gigii_Gii



luigi-seminara