





### Contrastive Sequential-Diffusion Learning: Non-linear and Multi-Scene Instructional Video Synthesis

WACV 2025
TUCSON, ARIZONA • FEB 28 - MAR 4



Vasco Ramos<sup>1</sup>, Yonatan Bitton<sup>2</sup>, Michal Yarom<sup>2</sup>, Idan Szpektor<sup>2</sup>, Joao Magalhaes<sup>1</sup>

<sup>1</sup>NOVA LINCS <sup>2</sup>Google Research

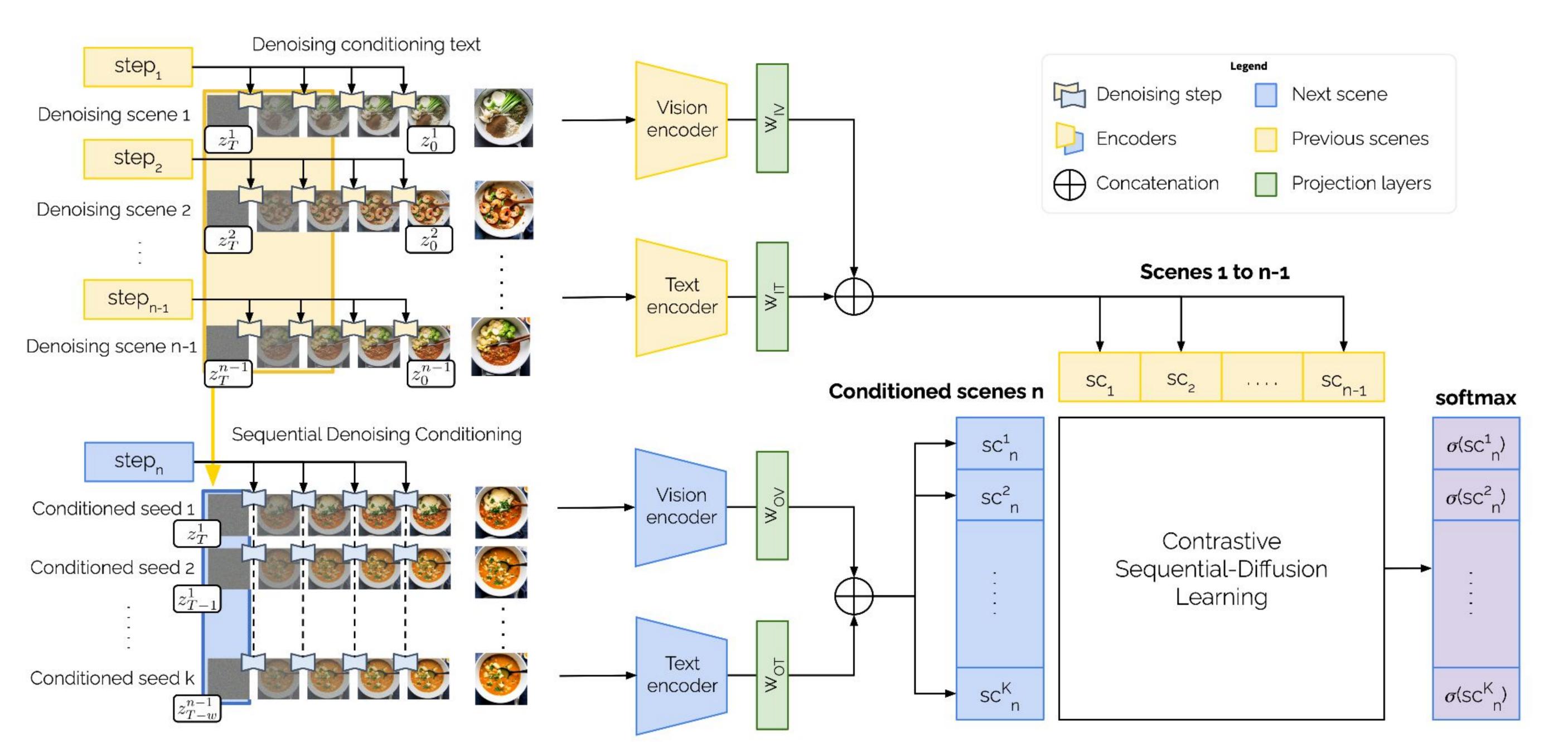
## Coherent Generation $\mathcal{L}_{CoSeD} = \mathbb{E}_{\mathbb{Z}_T^n, s_n, \varepsilon, t = T} \left[ \| \varepsilon - \varepsilon_\theta (\mathbb{Z}_T^n \mathbb{C}_n) \|_2^2 \right]$ Previous Latent Visual Caption Spices and vegetables in a pan Jamaican Callaloo and tomato on a pan

### Training

- Learn relationships between sequential steps across various tasks simultaneously.
- Ground truth information acts as both positive and negative pairs.
- Enable the model to leverage information from multiple tasks for improved learning.

### Task 1 s. V. P P N N N N N Task 2 scene 3 S. V. N N N N P P P N N N Scene 3

### Architecture



# COSED



Non-Linear

Linear

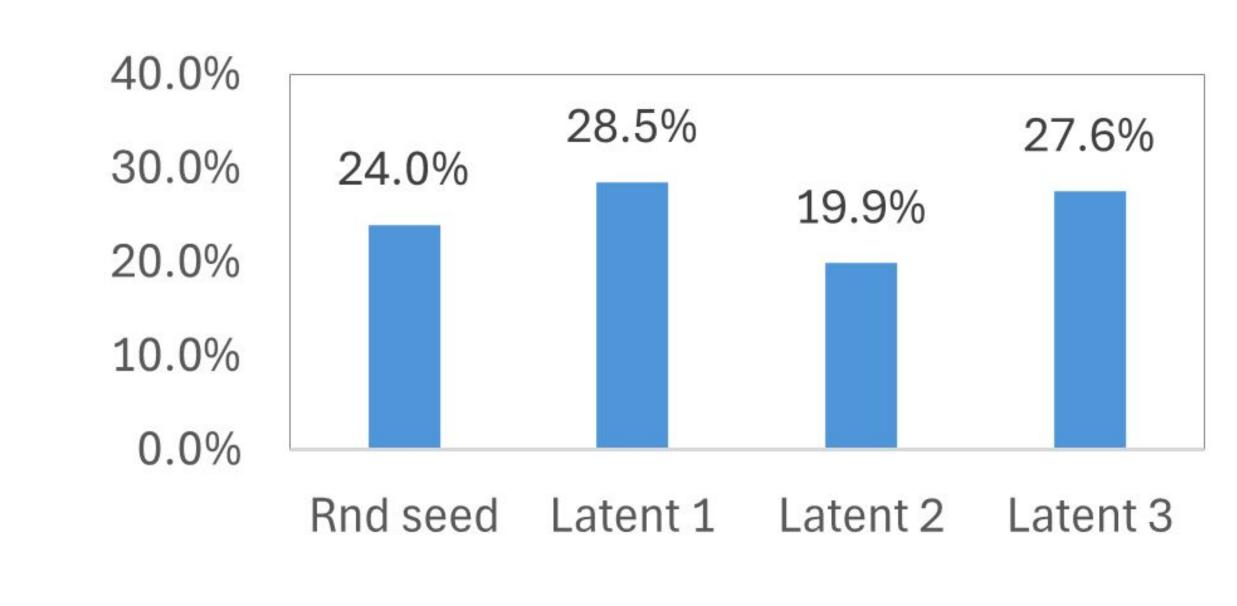
Specificalized in cooking recipes and DIY task videos.

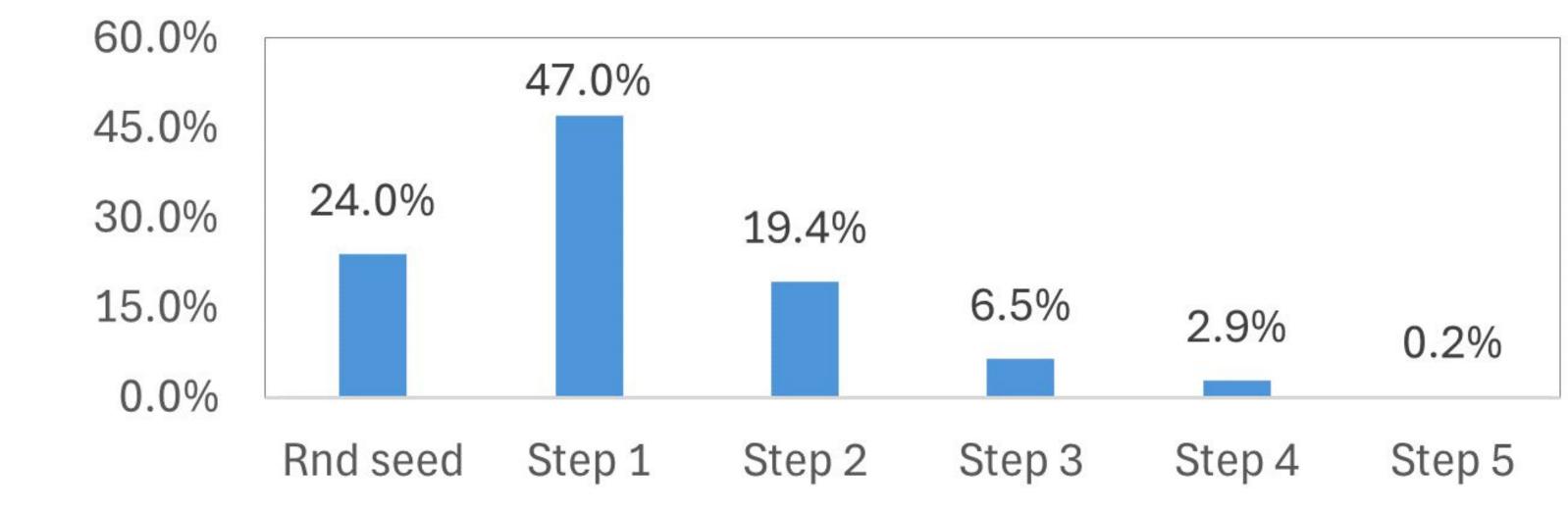
Contrastive selection ensures more coherent sequence generation.

Improved Seed Selection Method for better Visual and Semantic Coherence.

Achieves up to 90% preference in human evaluation compared to baselines.

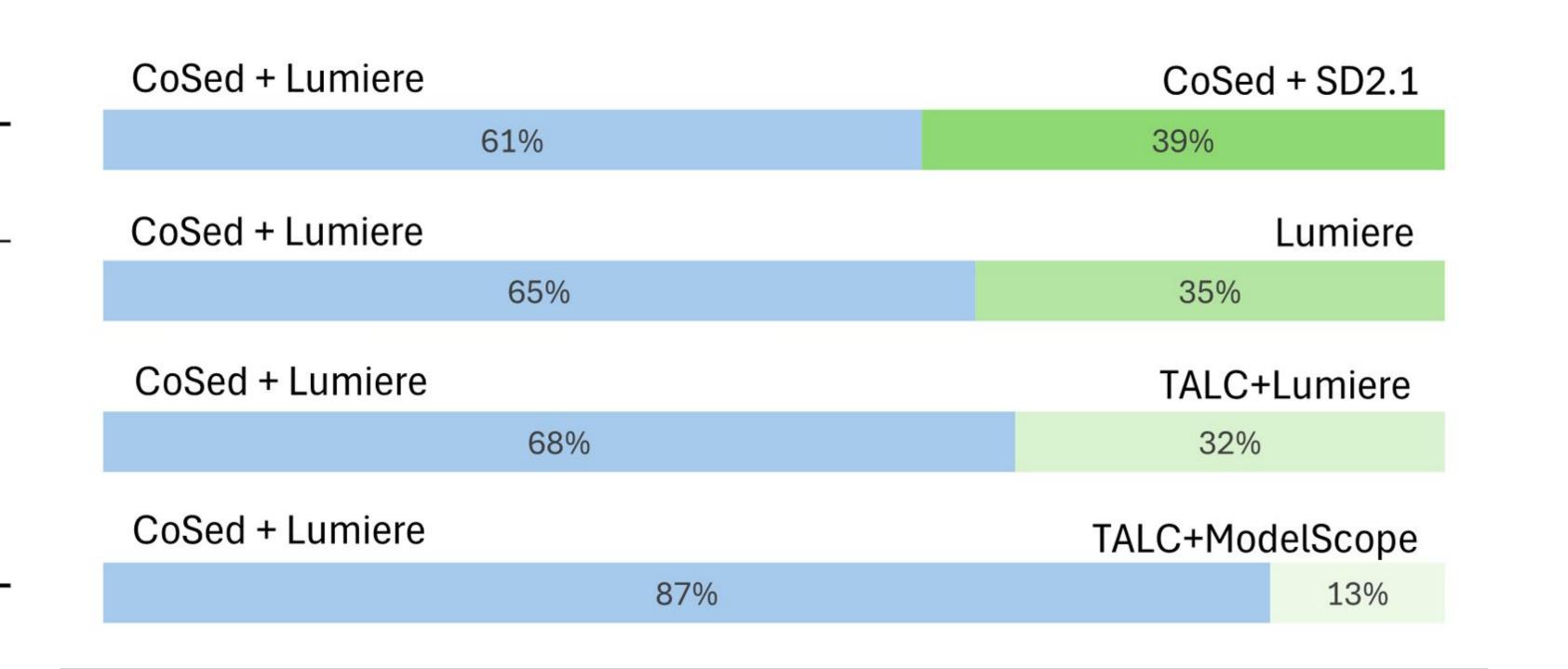
### Non-Linear Behaviour





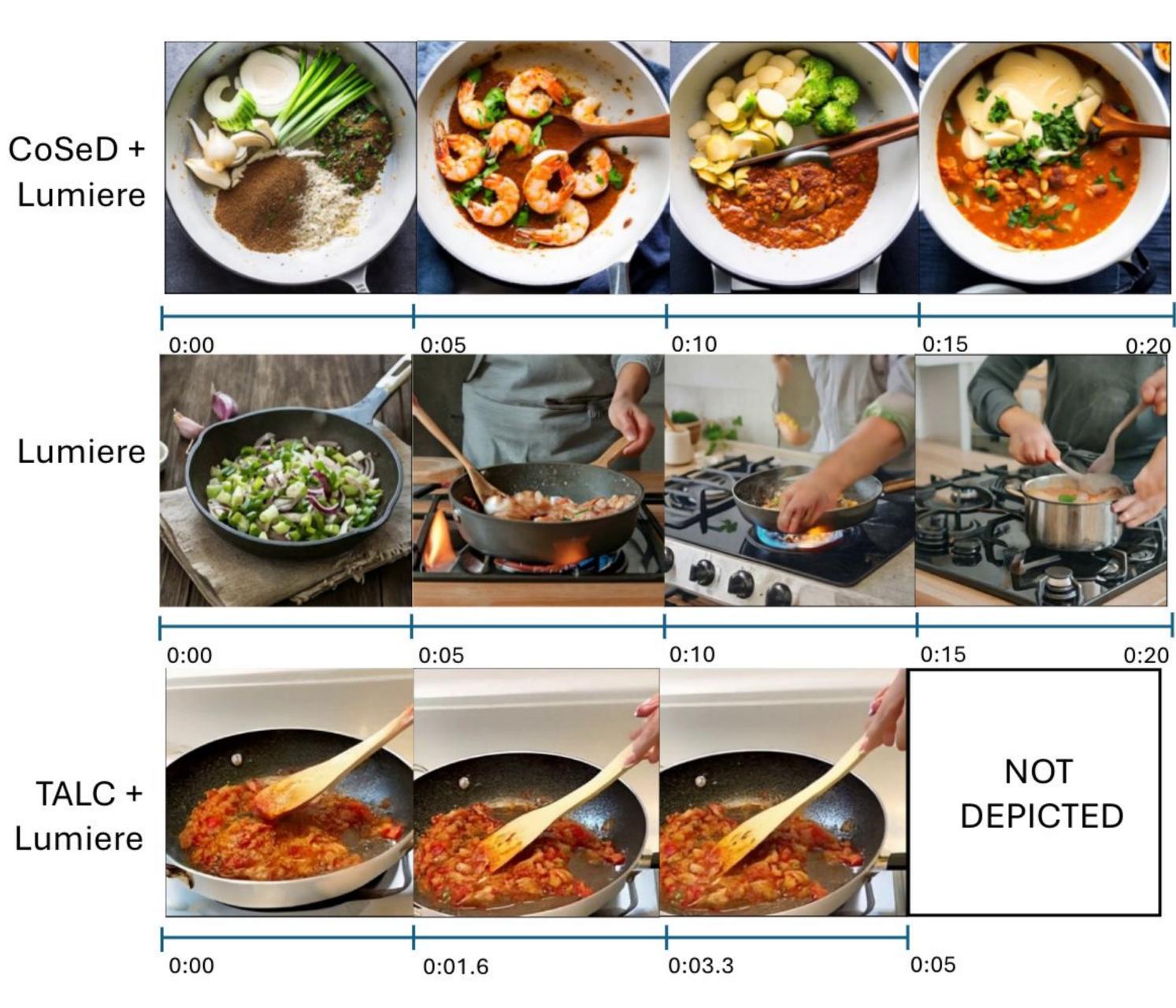
### Human Evaluation

| Methods           | Video<br>Length | Semantic<br>Consist. | Sequence<br>Consist. |
|-------------------|-----------------|----------------------|----------------------|
| CoSeD + Lumiere   | 20.8 s          | 85.0                 | 74.2                 |
| CoSeD + SVD       | 14.9 s          | 78.3                 | 69.2                 |
| TALC + ModelScope | 7.4 s           | 38.3                 | 50.8                 |
| TALC + Lumiere    | 5.0 s           | 30.0                 | 50.8                 |
| SD + SVD          | 14.9 s          | 80.0                 | 66.3                 |
| Lumiere           | 20.8 s          | <u>81.7</u>          | <u>72.9</u>          |



### Qualitative Results

- 1. Heat the Coconut oil, then add the Onion, Garlic and Scallion.
- 2. Add the Shrimp, and cook.
- 3. Turn the heat up and add the Jamaican Callaloo and Tomato.
- 4. After 10 minutes, taste for salt.



This work was partially supported by a Google Research Gift and by the FCT project NOVA LINCS Ref. (UIDB/04516/2020).