

# How Confident are Video Models?

## Empowering Video Models to Express their Uncertainty

Zhitong Mei<sup>1\*</sup>, Ola Shorinwa<sup>1\*</sup>, Anirudha Majumdar<sup>1</sup>

<sup>1</sup>Princeton University

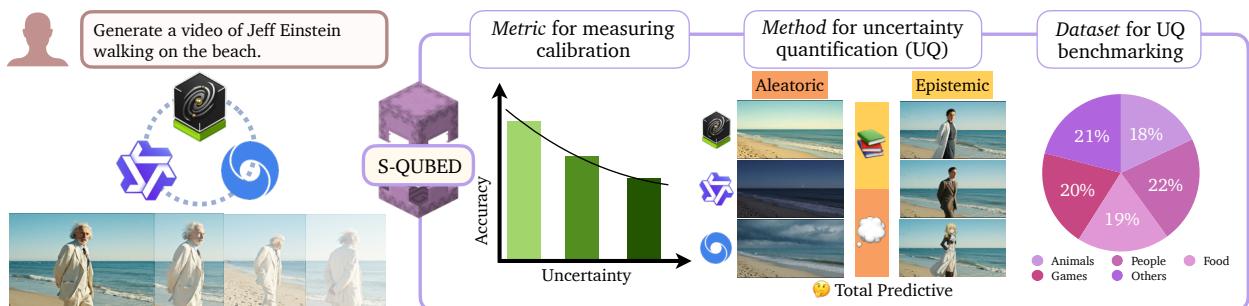
\*Equal contribution.

Generative video models demonstrate impressive text-to-video capabilities, spurring widespread adoption in many real-world applications. However, like large language models (LLMs), video generation models tend to *hallucinate*, producing plausible videos even when they are factually wrong. Although uncertainty quantification (UQ) of LLMs has been extensively studied in prior work, no UQ method for video models exists, raising critical safety concerns. To our knowledge, this paper represents the first work towards quantifying the uncertainty of video models. We present a framework for uncertainty quantification of generative video models, consisting of: (i) a metric for evaluating the calibration of video models based on robust rank correlation estimation with no stringent modeling assumptions; (ii) a black-box UQ method for video models (termed **S-QUBED**), which leverages latent modeling to rigorously decompose predictive uncertainty into its aleatoric and epistemic components; and (iii) a UQ dataset to facilitate benchmarking calibration in video models. By conditioning the generation task in the latent space, we disentangle uncertainty arising due to vague task specifications from that arising from lack of knowledge. Through extensive experiments on benchmark video datasets, we demonstrate that S-QUBED computes calibrated total uncertainty estimates that are negatively correlated with the task accuracy and effectively computes the aleatoric and epistemic constituents.

**Keywords:** Video Models, Uncertainty Quantification, Trustworthy Generative Models.

**Website:** [s-qubed.github.io](https://s-qubed.github.io)

**Code:** [github.com/irom-princeton/s-qubed](https://github.com/irom-princeton/s-qubed)



**Figure 1** Video models are unable to express their uncertainty, posing a critical limitation especially in tasks where they lack requisite knowledge. Here, the video model generates an inaccurate video (showing Albert Einstein), when prompted to generate a video of Jeff Einstein. To this end, we introduce a *metric* for evaluating the calibration of video models, a *calibrated uncertainty quantification method* (S-QUBED) which uses latent modeling to disentangle aleatoric and epistemic uncertainty, and a *UQ dataset* for benchmarking calibration.

# 1 Introduction

Recent advances in video generation models have led to huge strides in their capabilities [9, 27]. However, current text-to-video models tend to hallucinate, generating videos misaligned with the user intention, or disobeying physical laws. Despite this important limitation, existing video models are unable to express their own uncertainties, unlike LLMs, posing a crucial safety concern. We illustrate hallucinations in video models in Figure 1. When prompted to generate a video of Jeff Einstein walking on a beach, the video model generates a video of Albert Einstein, an entirely different person, without expressing any doubt in its output. We aim to address this critical challenge by empowering video models to express their uncertainty.

Specifically, we propose a framework for uncertainty quantification of video models, consisting of three fundamental components: First, we introduce a *metric* for evaluating the calibration of video models that directly assesses the alignment of the uncertainty estimates with the accuracy of the video generation task. Our metric estimates the rank correlation between uncertainty and task accuracy to measure the calibration error.

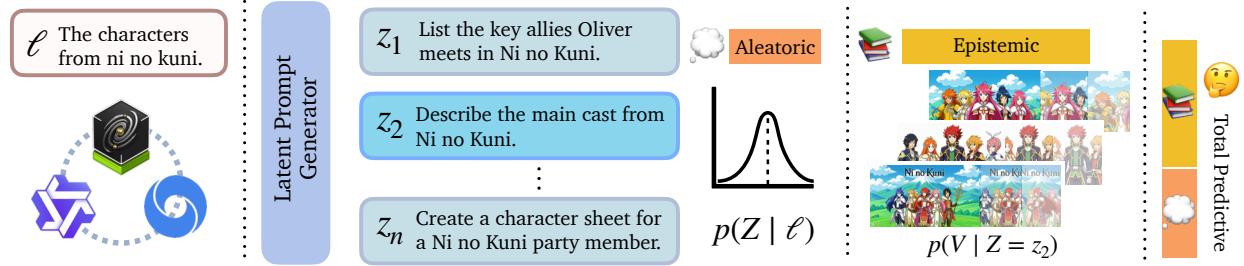
Second, we derive *S-QUBED* (Semantically-Quantifying Uncertainty with Bayesian Entropy Decomposition), a black-box uncertainty quantification method for video generation models, preserving amenability to the ever-increasing set of closed-source video models. Our key insight is to quantify uncertainty with latent modeling, enabling the rigorous decomposition of predictive uncertainty into its aleatoric and epistemic components. By mapping the input text prompt to a latent space, S-QUBED effectively distinguishes between uncertainty arising from ambiguous prompts and uncertainty arising from the model’s lack of knowledge. We demonstrate the calibration of S-QUBED’s estimates across a broad variety of video generation tasks.

Third, we curate a *UQ dataset* of about 40K videos across diverse tasks to facilitate benchmarking UQ methods for video models. We generate the data using the open-source model Cosmos-Predict2 [27]. We hope that the dataset drives research on uncertainty quantification of video models.

## 2 Related Work

**Uncertainty Quantification in Deep Learning.** Deep neural networks (DNNs) are generally difficult to interpret [20], motivating the development of UQ methods to examine the trustworthiness of their predictions [1]. UQ methods in deep learning can be broadly categorized into: *training-free* and *training-based* methods, which constitute a majority of existing work. Training-free methods estimate uncertainty without modifying the model’s architecture, training algorithm, or dataset, e.g., via perturbation techniques [23], dropout injection [19, 24], and test-time data augmentation [2, 38]. In contrast, training-based methods impose specific architectural design choices to enable uncertainty quantification using Bayesian Neural Networks (BNN) and can be further classified into three categories: (i) variational inference, (ii) Monte-Carlo Dropout, and (iii) Deep Ensemble methods. Assuming that the parameters (weights) of learned models are random variables, BNN methods [] apply Bayes’ rule to estimate a posterior distribution over these parameters given a prior distribution. However, the exact application of Bayes’ rule is typically intractable, giving rise to approximation techniques, e.g., variational inference [39], which approximates the posterior distribution using a parametric distribution; Monte-Carlo Dropout [12], which samples from the posterior distribution by zeroing-out some weights; and Deep Ensembles [18], which train multiple independent models to represent the posterior distribution. Despite their success, traditional UQ methods in deep learning are computationally expensive, limiting their applications in large generative models, e.g., large language models (LLMs) and vision-language models (VLMs). UQ methods for LLMs/VLMs generally leverage internal activations of these models, or utilize similarity-based metrics or natural-language inference techniques for more efficient UQ (see [32] for a detailed discussion).

**Uncertainty Quantification in Generative Image/Video Models.** Unlike DNNs and LLMs, UQ of generative image/video models has been relatively underexplored [11]. Prior work [7] extends Bayesian UQ techniques to denoising diffusion probabilistic models (DDPMs) in generative image modeling by learning a distribution of weights for the diffusion model, enabling the estimation of epistemic uncertainty through the variance across the model’s predictions. Similarly, other approaches [5] train latent diffusion models (diffusion ensembles)



**Figure 2 S-QUBED architecture.** Given a text prompt  $\ell$ , our goal is to quantify the uncertainty of the video generation model. We first generate  $n$  latent prompts consistent with  $\ell$  in line with the prompt refinement used by video models, modeling the aleatoric uncertainty as the entropy of the distribution over latent prompts. Then, for each latent prompt, we generate  $m$  videos, modeling the epistemic uncertainty as the conditional entropy of the distribution over generated videos. Finally, aggregating the two types of uncertainties yields the total predictive uncertainty.

for UQ by estimating the mutual information over a distribution of the models’ weights, analogous to deep ensembles. However, these training-based UQ methods are challenging to implement, given that diffusion models often have billions of parameters, creating significant computation overhead during training or inference. Drawing insights from black-box UQ methods for LLMs [4, 22, 25] which utilize similarity-based techniques for efficient UQ, PUNC [11] explores uncertainty quantification of generative image models in language space. By mapping generated images into language form using a VLM, PUNC leverages widely-used text-based similarity metrics [21, 41] to estimate epistemic and aleatoric uncertainty of text-to-image models. Although PUNC addresses the computation limitations of prior UQ methods for generative image models, PUNC is not applicable to video modes. To our knowledge, this work is the first exploration of UQ for video world models.

### 3 Problem Formulation

We examine uncertainty quantification of black-box text-conditioned video generation models, which map a text prompt  $\ell \in \mathcal{O}$  to a video  $v \in \mathcal{V}$  via an unknown stochastic model  $f_\theta : \mathcal{T} \mapsto \mathcal{V}$  parametrized by weights  $\theta$ . Specifically, the video generation process is described by the model:

$$v \sim f_\theta(V | \ell), \quad (1)$$

where  $v$  is sampled from the conditional distribution  $f_\theta$ . For an input prompt  $v$ , the video generation model has a measure of doubt (uncertainty) associated with the sampled video output  $v$ . This uncertainty arises from a variety of sources, e.g., vagueness in the conditioning input  $\ell$ , randomness in the physical evolution of the real-world, limited training data, etc. In this work, we are interested in quantifying the *total* predictive uncertainty associated with  $v$ , which can be broadly classified into two categories: *aleatoric* uncertainty and *epistemic* uncertainty.

### 4 Uncertainty Quantification of Generative Video Models

We present S-QUBED, an efficient method for uncertainty quantification of video generation models, summarized in Figure 2. Without loss of generality, we can decompose the video generation model in Equation (1) using a latent variable  $z \in \mathcal{Z}$ , modeling the video generation as a two-step process. In the first step,  $z$  is sampled from the probability distribution  $p(Z | \ell)$  conditioned on the input prompt  $\ell$ . In the second step, the video model samples the output video  $v$  from the probability distribution  $p(V | Z = z)$ . Note that the application of latent variables is standard in generative modeling, e.g., in variational Bayesian learning [6, 17, 33], enabling efficient learning and analysis of complex data-generation distributions. Consequently, we can rewrite Equation (1) in the form:

$$f_\theta(V | \ell) = \int_{z \in \mathcal{Z}} p(V | z, \ell) p(z | \ell) dz = \int_{z \in \mathcal{Z}} p(V | z) p(z | \ell) dz, \quad (2)$$

where we assumed conditional independence of  $V$  and  $\ell$ , given  $z$ .

Note that the video generation model described by [Equation \(2\)](#) is not limiting. In fact, state-of-the-art text-to-video models refine a user’s prompt using an LLM to generate a much more detailed prompt that is passed into the video generation model. Hence, we can interpret [Equation \(2\)](#) as first sampling an instance of a fully-specified prompt  $z$  from the conditional distribution defined by the input prompt  $\ell$ , e.g., given the input prompt “a cat doing something,”  $z$  may be the more specific prompt “a cat licking its paws before turning to the camera and meowing...” Subsequently, the video model generates the output video conditioned on  $z$ .

**Proposition 1** (Uncertainty Decomposition). *Define the total predictive uncertainty in the output video as the differential entropy  $h(V | \ell)$  of the distribution  $f_\theta(V | \ell)$ . Then, this quantity can be decomposed as:*

$$h(V | \ell) = h(V | Z) + h(Z | \ell), \quad (3)$$

where  $h(V | Z)$  represents the epistemic uncertainty in  $v$ , and  $h(Z | \ell)$  the aleatoric uncertainty.

This is a standard decomposition. We provide the proof in [Appendix A.2](#) for completeness. In the rest of this section, we introduce our approach to estimating these components.

## 4.1 Aleatoric Uncertainty

Aleatoric uncertainty encompasses irreducible randomness from the vagueness (lack of sufficient specificity) of the conditioning inputs, e.g., “generate a video of a cat doing something.” In video generation, vagueness in the input prompt increases the randomness of the conditional probability distribution  $p(Z | \ell)$ , which is represented by the second term  $h(Z | \ell)$  in [Equation \(3\)](#). Note that  $h(Z | \ell)$  is independent of  $v$  since the source of uncertainty arises from the input prompt independent of the second stage of the video generation, e.g., the denoising process in video diffusion models. In particular, randomness in  $Z$  cannot be reduced by training the video model on additional data under the assumption that we can model  $p(Z | \ell)$  *almost* exactly.

As a measure of aleatoric uncertainty, we would expect  $h(Z | \ell)$  to be positively correlated with the vagueness of the input prompt. For example, consider two input prompts:  $\ell_1$  = “a cat napping” and  $\ell_2$  = “a cat doing something”. With  $\ell_1$ , the pdf of  $p(Z | \ell_1)$  will be concentrated on the set:

$$\mathcal{A}(\ell_1) = \{“a black cat napping”, “a cat napping on a couch”, “a cat snoring on a couch”, …\}. \quad (4)$$

However, with  $\ell_2$ , the pdf of  $p(Z | \ell_2)$  will be concentrated on the set:

$$\mathcal{A}(\ell_2) = \{“a black cat jumping”, “a cat eating on a couch”, “a cat meowing next to a door”, …\}. \quad (5)$$

Note that the elements of  $\mathcal{A}(\ell_1)$  are more semantically-related (since  $\ell_1$  is more specific) and are thus closer in the language (semantic) embedding space compared to elements in  $\mathcal{A}(\ell_2)$ . Hence,  $p(Z | \ell_1)$  will have a lower entropy relative to  $p(Z | \ell_2)$ .

**Modeling the conditional latent distribution.** To compute  $h(Z | \ell)$ , we need to define a class of probability distributions that describe the latent-generation process. In this work, we model  $p(Z | \ell)$  in a language embedding space using the Von-Mises Fisher (VMF) distribution [10, 15], drawing insights from prior work [3, 14, 31].

The Von-Mises Fisher (VMF) distribution describes a  $n$ -dimensional probability distribution on the  $(n - 1)$ -sphere over unit vectors embedded in  $\mathbb{R}^n$ , with the probability density function (pdf):

$$f_n(x, \mu, \kappa) = C_n(\kappa) \exp(\kappa \mu^\top x), \quad (6)$$

with parameters  $\mu$  and  $\kappa$  denoting the mean direction and concentration parameters, where:

$$C_n(\kappa) = \frac{\kappa^{n/2-1}}{(2\pi)^{n/2} I_{n/2-1}(\kappa)}, \quad (7)$$

with  $I_{n/2-1}$  representing the modified Bessel function of the first kind. The concentration parameter functions analogously to the inverse variance, providing a measure of the spread of the distribution.

We need samples from  $p(Z | \ell)$  to fit the VMF distribution. Collecting such data is typically prohibitively expensive. To overcome this challenge, we leverage LLMs as cost-effective generative models of  $p(Z | \ell)$ , noting that video models generally use LLMs to refine prompts prior to generating videos.

Specifically, given an input prompt  $\ell$ , we generate  $N$  *compatible*-but-more-specific prompts from an LLM. A generated prompt is *compatible* with the input prompt if the generated prompt is consistent with, i.e., *entails*, the input prompt. However, the converse need not be true: the input prompt might be underspecified. Subsequently, we compute language embeddings from an embedding model, e.g., SentenceFormer [30]. Although we could directly fit a VMF to the language embeddings, we project the language embeddings to a lower-dimensional subspace  $\mathbb{R}^n$  using principal component analysis (PCA) to avoid numerical instability associated with high-dimensional spaces. We estimate the parameters  $\mu$  and  $\kappa$  of the VMF distribution in closed-form using approximate methods [15, 35], circumventing iterative optimization methods.

**Estimating Aleatoric Uncertainty.** Given  $p(Z | \ell)$ , we can compute the aleatoric uncertainty  $h(Z | \ell)$  of  $v$  in closed-form via:

$$h(Z | \ell) = -\log(C_n(\kappa)) - \frac{\kappa}{\mu_{z|\ell}} \mathbb{E}_Z[Z | \ell](\kappa), \quad (8)$$

where  $Z \sim \text{VMF}(\mu, \kappa)$  and  $C_n$  represents the normalization constant given by Equation (7). The expected value of the VMF is given by  $\mathbb{E}_Z[Z | \ell](\kappa) = W_n(\kappa)\mu_{z|\ell}$ , where  $W_n = \frac{I_{n/2}(\kappa)}{I_{n/2-1}(\kappa)}$  with the modified Bessel function of the first kind  $I_{n/2}$ . We summarize the method for computing aleatoric uncertainty in Algorithm 1.

---

**Algorithm 1:** S-QUBED: Aleatoric Uncertainty Quantification of Generative Video Models

---

**AleatoricUncertainty** ( $f, \ell$ ):

<b>Input:</b> Video Model $f$ , Input Prompt $\ell$ ;	
<b>Output:</b> Aleatoric Uncertainty $h(Z   \ell)$ ;	
$\mathcal{A}(\ell) \leftarrow \text{Embed(LLM}(\ell))$ ;	<i>// Construct <math>\mathcal{A}(\ell)</math> from an LLM/VLM</i>
$\mu_{z \ell}, \kappa_{z \ell} \leftarrow \text{VFit}(\mathcal{A}(\ell))$ ;	<i>// Estimate <math>p(Z   \ell)</math> with a VMF</i>
$h(Z   \ell) \leftarrow \text{Equation (8)}$ ;	<i>// Compute aleatoric uncertainty <math>h(Z   \ell)</math></i>
<b>return</b> $h(Z   \ell)$ ;	

---

## 4.2 Epistemic Uncertainty

Epistemic uncertainty represents the measure of doubt associated with a lack of knowledge, which generally results from insufficient training data (e.g., Figure 1). As a result, epistemic uncertainty is *reducible* by providing additional training data to the model. In Equation (3),  $h(V | Z)$  represents the epistemic uncertainty of the generated video  $v$ , where the uncertainty arises from the limited knowledge of the video model about concepts expressed by the latent variable  $z \in \mathcal{Z}$ .

For example, consider a video model trained entirely on internet videos of cats and dogs performing different activities, e.g., running, eating, jumping, meowing/barking. Now, when asked to generate a video of “a lion roaring in the wild”, the video model might generate different videos across different runs, with some showing a large cat meowing in a park with significant tree canopy, others showing a cat making *barking-like* sounds in a forest, etc. Although the generated videos are all conditioned on semantically-consistent latent variables, the generated videos might be semantically-inconsistent, since the video model has not been trained on videos of lions. This uncertainty in the generated videos can be described as *epistemic* and is captured by the entropy term  $h(V | Z)$ .

**Estimating Epistemic Uncertainty.** Note that we can express  $h(V | Z)$  in the form:

$$h(V | Z) = \mathbb{E}_{z \sim p(z|\ell)}[h(V | Z = z)], \quad (9)$$

which can be interpreted as the expected entropy of the distribution of generated videos conditioned on sampled latent states  $z$  from the conditional distribution  $p(z | \ell)$ . Computing  $h(V | Z)$  is challenging for two

reasons: (i) we do not have an explicit model of  $p(V | Z = z)$  which is required to compute  $h(V | Z = z)$ , and (ii) even with an analytical expression for  $p(V | Z = z)$ , computing  $h(V | Z)$  would require evaluating a double integral, which is intractable in general.

To address the first challenge, we approximate the conditional distribution  $p(V | Z = z)$  using a VMF distribution with the parameters  $\mu$  and  $\kappa$  estimated from samples drawn from the video model. Likewise, we approximate the expectation in [Equation \(9\)](#) using Monte-Carlo sampling to address the second challenge, which we describe in greater detail.

First, we sample a set of latent variables  $\mathcal{E}_{z|\ell}$  conditioned on the input prompt  $\ell$  from the distribution  $p(Z | \ell)$ , with each  $z \in \mathcal{E}_{z|\ell}$  representing specific instances of prompts entailing the input prompt. For each  $z$ , we estimate the distribution  $p(V | Z = z)$  by generating a set of videos  $\mathcal{E}_{v|z}$  from the video model, conditioned on  $z$ . Subsequently, we embed these videos with a video embedding model, e.g., S3D [26] and fit a VMF to the samples in  $\mathcal{E}_{v|z}$ . Afterwards, we compute the entropy  $h(V | Z = z)$  with:

$$h(V | z) = -\log(C_n(\kappa_{v|z})) - \frac{\kappa_{v|z}}{\mu_{v|z}} \mathbb{E}_{v|z}[V | Z = z](\kappa_{v|z}), \quad (10)$$

using the estimated VMF parameters  $\mu_{v|z}$  and  $\kappa_{v|z}$ . Finally, we compute an empirical estimate of the expectation of  $h(V | Z = z)$  over  $z$  sampled from  $p(Z | \ell)$ . We outline these steps in [Algorithm 2](#).

---

**Algorithm 2:** S-QUBED: Epistemic Uncertainty Quantification of Generative Video Models

---

**EpistemicUncertainty** ( $f, \ell$ ):

```

Input: Video Model  $f$ , Input Prompt  $\ell$ ;
Output: Epistemic Uncertainty  $h(V | z)$ ;
 $\mathcal{E}_{z|\ell} \leftarrow \text{Embed(LLM}(\ell)\text{)}; \quad // \text{Construct } \mathcal{E}_{z|\ell} \text{ from an LLM/VLM}
\textbf{foreach } z \in \mathcal{E}_{z|\ell} \textbf{ do}
    \mathcal{E}_{v|z} \leftarrow \text{Embed}(f(V | z)); \quad // \text{Construct } \mathcal{E}_{v|z} \text{ from } f
    \mu_{v|z}, \kappa_{v|z} \leftarrow \text{VFit}(\mathcal{E}_{v|z}); \quad // \text{Estimate } p(V | Z = z) \text{ from } f
    h(V | Z = z) \leftarrow \textcolor{red}{\text{Equation (10)}}; \quad // \text{Compute entropy } h(V | Z = z)
\textbf{end}
h(V | Z) \leftarrow \textcolor{red}{\text{Equation (9)}}; \quad // \text{Compute epistemic uncertainty } h(V | Z)
\textbf{return } h(V | Z);$ 
```

---

## 5 Experiments

We examine the effectiveness of S-QUBED in uncertainty quantification of generative video models, specifically exploring the following questions: (i) *How do we evaluate uncertainty calibration of video models?* (ii) *Are the total predictive uncertainty estimates computed by S-QUBED calibrated?* (iii) *Can S-QUBED effectively estimate both aleatoric and epistemic uncertainty?*

### 5.1 Evaluation Setup

We describe the datasets, models, and metrics used in evaluating our proposed method.

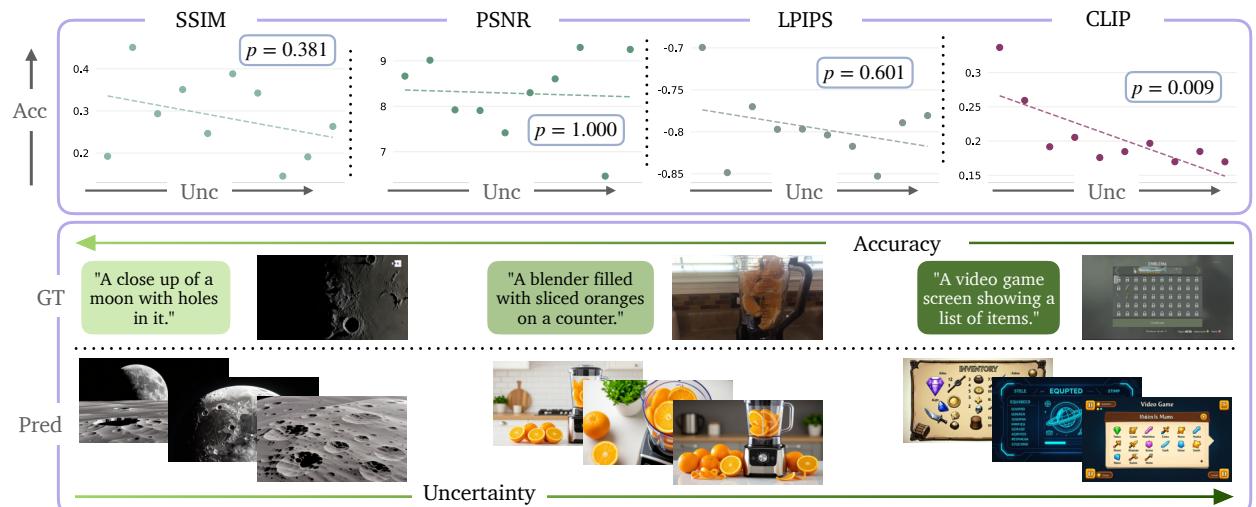
**Datasets.** We evaluate S-QUBED on two large-scale video generation datasets, VidGen-1M [36] and Panda-70M [8]. Using GPT-5-nano [29], we classify the videos in each dataset into five broad categories: animals, food, games, people, and other, a standard approach with video datasets. We subsample about 200 video generation tasks uniformly from each category for evaluation. To address issues with missing video data/metadata in some of the datasets, we sample additional videos from other categories, minimally changing the uniformity of the evaluation dataset.

**Implementation.** We evaluate S-QUBED on the Cosmos-Predict2 video model [28] using the official implementation, which utilizes a text-to-image-to-video pipeline for text conditioning that generates an image from a text prompt, which is used as input to an image-to-video model. Although we explored alternative generative

video models, e.g., Veo3 [9], none were compatible with our experiments, either due to limitations on the number of permissible generation requests or prohibitive compute cost. For example, Veo3 only supports generating between 5-20 videos per month, which is insufficient for the scale of our evaluations. We implement our proposed method by sampling 10 latent states,  $z_{1:10} \sim p(Z|\ell)$ , and subsequently 10 generated videos per latent state,  $v_{1:10}^i \sim p(v_{1:10}^i|Z = z_i)$ .

## 5.2 How do we evaluate uncertainty calibration of video models?

Uncertainty calibration of video generation models has been underexplored, evidenced by the lack of purpose-specific calibration metrics. Widely-used calibration metrics, such as the expected calibration error (ECE) and maximum calibration metrics (MCE) apply only to evaluation settings with discrete ground-truth answers and errors, e.g., with multiple-choice questions, making them unsuitable in video generation tasks with real-valued task errors. Consequently, we propose appropriate metrics for evaluating the calibration of the uncertainty estimates of video models. Specifically, we examine the Kendall rank correlation (Kendall’s  $\tau$ ) [16] between the video model’s uncertainty estimates and an applicable accuracy metric, which captures the degree of monotonicity between uncertainty and accuracy. We do not utilize Pearson’s rank correlation [13] due to its assumptions of linearity and normally-distributed data and likewise do not use the Spearman’s rank correlation coefficient [34] due to its high sensitivity to outliers.



**Figure 3 Calibration Metrics for Video Models.** *Top:* We examine the statistical significance of the Kendall rank correlation between uncertainty and widely-used perceptual metrics. We find that the CLIP cosine similarity score provides the most significant correlation. *Bottom:* With the CLIP accuracy metric, we observe that low human-annotated uncertainty corresponds to smaller variance in the generated videos and greater accuracy with respect to the ground-truth video. As uncertainty increases, video prediction accuracy decreases.

To compute the rank correlation coefficient, we use the SSIM, PSNR, LPIPS, and CLIP score metrics. To identify the best metric for assessing calibration, we select 10 generation tasks from the Panda-70M datasets and rank the tasks in order of increasing uncertainty based on the vagueness of the text prompt for the task. Note that the vagueness in the prompt directly corresponds to aleatoric uncertainty, making it an effective proxy measure. Given the human-annotated rankings, we compute the Kendall rank correlation between uncertainty and each accuracy metric along with a  $p$ -value, which provides a measure of the statistical significance of the correlation. While Panda-70M dataset consists of tasks with a broad range of descriptive detail from vague to very specific, VidGen-1M consists of relatively well-detailed tasks. As a result, we do not sample from VidGen-1M, given the less observable variation in the aleatoric uncertainty. We sample the tasks from Panda-70M dataset to retain the distribution of instruction detail.

We summarize our results in Figure 3. In Panda-70M, the CLIP score metric is strongly negatively correlated with uncertainty at the 99% significance level. In contrast, the other perceptual metrics lack a statistically significant correlation with uncertainty. This finding is not entirely surprising, since CLIP captures semantic

information that better reflects the accuracy of the generation task, unlike the other perceptual metrics which are more susceptible to differences in visual changes.

Moreover, we visualize the text prompt, ground-truth video, and the first frame of the generated videos for a few tasks in [Figure 3](#), ranging from low to high uncertainty (rank). We observe that when uncertainty is low, the model tends to generate very similar videos, which are also close to the ground-truth, resulting in high accuracy with respect to the CLIP score. As we vary the uncertainty of the model, we observe greater variance in the generated videos accompanied by notably lower CLIP scores (compared to the other metrics), further demonstrating the utility of the CLIP score as an accuracy metric.

### 5.3 Are our uncertainty estimates calibrated?

We examine the calibration of our uncertainty estimates in VidGen-1M and Panda-70M, using the CLIP score accuracy metric given its effectiveness in assessing calibration. We first compute the total predictive uncertainty associated with each video task using S-QUBED, and then evaluate the Kendall rank correlation. We define the accuracy of each task as the mean CLIP score across all generated videos for that task.

[Figure 4](#) (left) presents results for Panda-70M. We observe a statistically significant negative correlation (99% confidence level) between the total uncertainty computed using S-QUBED and the CLIP score, demonstrating calibration of the uncertainty estimates. The results highlight that as the uncertainty of the video model decreases, its accuracy increases. Likewise, in VidGen-1M, the total predictive uncertainty is negatively correlated with the CLIP score at the 89.9% confidence level. From [Figure 4](#), we see that when the total predictive uncertainty estimates is small (“A”), the video model generates more accurate videos; in contrast, in tasks with high estimated uncertainty (“B”), the video model is less accurate.



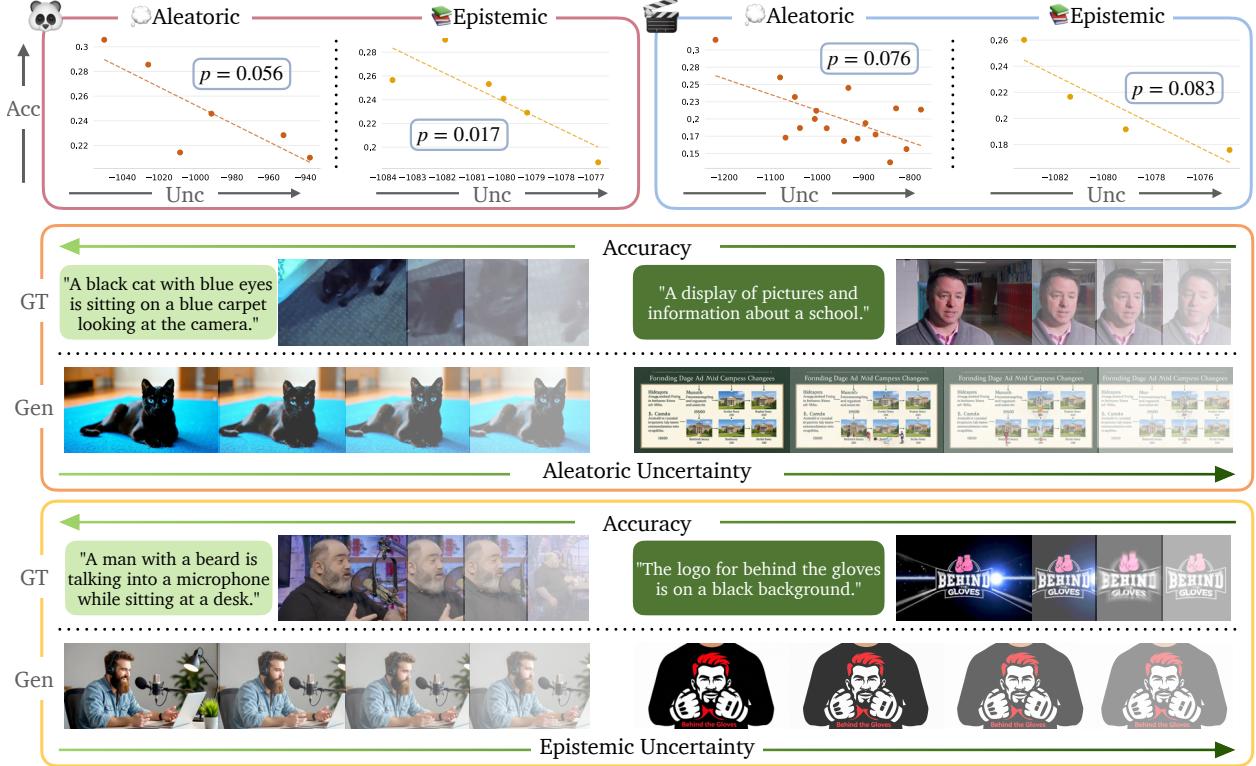
**Figure 4 Total Predictive Uncertainty for Video Models.** We assess the calibration of the total predictive uncertainty computed by S-QUBED. *Top:* correlation between video prediction accuracy and total uncertainty for Panda-70M and VidGen-1M . We observe a statistically significant correlation between accuracy and uncertainty for both datasets, signified by the small  $p$ -values. *Bottom:* visualization of two samples from Panda-70M.

### 5.4 Can S-QUBED effectively estimate both aleatoric and epistemic uncertainty?

We examine the performance of S-QUBED in decomposing total uncertainty into aleatoric and epistemic uncertainty. To effectively assess calibration of aleatoric uncertainty, we consider a subset of each dataset where the epistemic uncertainty is almost zero and compute the rank correlation between the aleatoric uncertainty of these samples and the CLIP score. Likewise, to evaluate calibration of epistemic uncertainty, we compute the rank correlation between the epistemic uncertainty and the CLIP score for samples with

relatively zero aleatoric uncertainty. In practice, we select samples with the lowest aleatoric or epistemic uncertainty, accordingly.

In [Figure 5](#), we visualize the Kendall rank correlation between the aleatoric and epistemic uncertainty and the CLIP score in both datasets. In Panda-70M, we find that aleatoric and epistemic uncertainty are negatively correlated with accuracy at the 94.5% and 98.3% confidence level. Similarly, in VidGen-1M, we observe a statistically significant negative correlation between aleatoric and epistemic uncertainty and the accuracy at the 92.3% and 91.7%, respectively. These results highlight that S-QUBED can decompose total uncertainty effectively into its aleatoric and epistemic components



**Figure 5 Disentangling Aleatoric and Epistemic Uncertainty for Video Models.** We demonstrate the calibration of the aleatoric uncertainty estimates of S-QUBED in tasks with no epistemic uncertainty, showing statistically significant negative correlation. We do the same for epistemic uncertainty.

Further, we visualize text prompts, ground-truth-videos, and generated videos in tasks with low and high estimated aleatoric uncertainty. We observe that in the low-uncertainty case, the video model achieves high accuracy, unlike the high-uncertainty case, where the prediction accuracy is significantly lower. Similarly, we provide some visualizations in the case with low and high estimated epistemic uncertainty, showing the negative correlation between S-QUBED’s estimated epistemic uncertainty and video prediction accuracy. Notably, the model does not know the specific prompt “Behind the Gloves” logo, unlike predicting the person in the human-centric videos.

## 6 Conclusion

We present a framework for empowering video models to express their uncertainty, a critical capability for safety. Concretely, we introduce a metric for measuring the calibration of UQ methods for video models and present a calibrated UQ method for video models. Our methods utilize latent modeling to estimate both aleatoric and epistemic uncertainty, without making any limiting assumptions. Further, we provide an open-source video dataset for benchmarking UQ methods for video models. Our experiments demonstrate the calibration of our proposed method and its effectiveness in disentangling aleatoric and epistemic uncertainty.

## 7 Limitations and Future Work

S-QUBED requires generating multiple videos from the video model to estimate epistemic uncertainty, which poses some computational overhead. Future work will explore more efficient strategies for sampling videos from the video model, e.g., in the latent space of the video model. Beyond the two benchmark datasets considered in this work, we will explore extensions to new datasets to augment the UQ dataset curated for benchmarking calibration. In addition, future work will examine the application of our method to new open-source models, as they become available.

## References

- [1] Moloud Abdar, Farhad Pourpanah, Sadiq Hussain, Dana Rezazadegan, Li Liu, Mohammad Ghavamzadeh, Paul Fieguth, Xiaochun Cao, Abbas Khosravi, U Rajendra Acharya, et al. A review of uncertainty quantification in deep learning: Techniques, applications and challenges. *Information fusion*, 76:243–297, 2021.
- [2] Murat Seckin Ayhan and Philipp Berens. Test-time data augmentation for estimation of heteroscedastic aleatoric uncertainty in deep neural networks. In *Medical Imaging with Deep Learning*, 2018.
- [3] Arindam Banerjee, Inderjit S Dhillon, Joydeep Ghosh, Suvrit Sra, and Greg Ridgeway. Clustering on the unit hypersphere using von mises-fisher distributions. *Journal of Machine Learning Research*, 6(9), 2005.
- [4] Evan Becker and Stefano Soatto. Cycles of thought: Measuring llm confidence through stable explanations. *arXiv preprint arXiv:2406.03441*, 2024.
- [5] Lucas Berry, Axel Brando, and David Meger. Shedding light on large generative networks: Estimating epistemic uncertainty in diffusion models. In *The 40th Conference on Uncertainty in Artificial Intelligence*, 2024.
- [6] Apratim Bhattacharyya, Michael Hanselmann, Mario Fritz, Bernt Schiele, and Christoph-Nikolas Straehle. Conditional flow variational autoencoders for structured sequence prediction. *arXiv preprint arXiv:1908.09008*, 2019.
- [7] Matthew Chan, Maria Molina, and Chris Metzler. Estimating epistemic and aleatoric uncertainty with a single model. *Advances in Neural Information Processing Systems*, 37:109845–109870, 2024.
- [8] Tsai-Shien Chen, Aliaksandr Siarohin, Willi Menapace, Ekaterina Deyneka, Hsiang-wei Chao, Byung Eun Jeon, Yuwei Fang, Hsin-Ying Lee, Jian Ren, Ming-Hsuan Yang, et al. Panda-70m: Captioning 70m videos with multiple cross-modality teachers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13320–13331, 2024.
- [9] DeepMind. Veo-3: A text-to-video generation system with audio. Technical Report Tech Report, DeepMind / Google, 2025. Accessed: YYYY-MM-DD.
- [10] Ronald Aylmer Fisher. Dispersion on a sphere. *Proceedings of the Royal Society of London. Series A. Mathematical and Physical Sciences*, 217(1130):295–305, 1953.
- [11] Gianni Franchi, Nacim Belkhir, Dat Nguyen Trong, Guoxuan Xia, and Andrea Pilzer. Towards understanding and quantifying uncertainty for text-to-image generation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 8062–8072, 2025.
- [12] Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059. PMLR, 2016.
- [13] Francis Galton. Note on regression and correlation. *Proceedings of the Royal Society of London*, 58:240–242, 1895.
- [14] Siddharth Gopal and Yiming Yang. Von mises-fisher clustering models. In *International Conference on Machine Learning*, pages 154–162. PMLR, 2014.
- [15] Peter Edmund Jupp and KV Mardia. A unified view of the theory of directional statistics, 1975-1988. *International Statistical Review/Revue Internationale de Statistique*, pages 261–294, 1989.
- [16] Maurice G Kendall. A new measure of rank correlation. *Biometrika*, 30(1-2):81–93, 1938.
- [17] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.

- [18] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. *Advances in neural information processing systems*, 30, 2017.
- [19] Emanuele Ledda, Giorgio Fumera, and Fabio Roli. Dropout injection at test time for post hoc uncertainty quantification in neural networks. *Information Sciences*, 645:119356, 2023.
- [20] Xuhong Li, Haoyi Xiong, Xingjian Li, Xuanyu Wu, Xiao Zhang, Ji Liu, Jiang Bian, and Dejing Dou. Interpretable deep learning: Interpretation, interpretability, trustworthiness, and beyond. *Knowledge and Information Systems*, 64(12):3197–3234, 2022.
- [21] Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81, 2004.
- [22] Zhen Lin, Shubhendu Trivedi, and Jimeng Sun. Generating with confidence: Uncertainty quantification for black-box large language models. *arXiv preprint arXiv:2305.19187*, 2023.
- [23] Yifei Liu, Rex Shen, and Xiaotong Shen. Novel uncertainty quantification through perturbation-assisted sample synthesis. *IEEE transactions on pattern analysis and machine intelligence*, 46(12):7813–7824, 2024.
- [24] Antonio Loquercio, Mattia Segu, and Davide Scaramuzza. A general framework for uncertainty estimation in deep learning. *IEEE Robotics and Automation Letters*, 5(2):3153–3160, 2020.
- [25] Potsawee Manakul, Adian Liusie, and Mark JF Gales. Selfcheckgpt: Zero-resource black-box hallucination detection for generative large language models. *arXiv preprint arXiv:2303.08896*, 2023.
- [26] Antoine Miech, Jean-Baptiste Alayrac, Lucas Smaira, Ivan Laptev, Josef Sivic, and Andrew Zisserman. End-to-end learning of visual representations from uncurated instructional videos. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9879–9889, 2020.
- [27] NVIDIA, :, Niket Agarwal, Arslan Ali, Maciej Bala, Yogesh Balaji, Erik Barker, Tiffany Cai, Prithvijit Chat-topadhyay, Yongxin Chen, Yin Cui, Yifan Ding, Daniel Dworakowski, Jiaoqiao Fan, Michele Fenzi, Francesco Ferroni, Sanja Fidler, Dieter Fox, Songwei Ge, Yunhao Ge, Jinwei Gu, Siddharth Gururani, Ethan He, Jiahui Huang, Jacob Huffman, Pooya Jannaty, Jingyi Jin, Seung Wook Kim, Gergely Klár, Grace Lam, Shiyi Lan, Laura Leal-Taixe, Anqi Li, Zhaoshuo Li, Chen-Hsuan Lin, Tsung-Yi Lin, Huan Ling, Ming-Yu Liu, Xian Liu, Alice Luo, Qianli Ma, Hanzi Mao, Kaichun Mo, Arsalan Mousavian, Seungjun Nah, Sriharsha Niverty, David Page, Despoina Paschalidou, Zeeshan Patel, Lindsey Pavao, Morteza Ramezanali, Fitsum Reda, Xiaowei Ren, Vasanth Rao Naik Sabavat, Ed Schmerling, Stella Shi, Bartosz Stefanik, Shitao Tang, Lyne Tchapmi, Przemek Tredak, Wei-Cheng Tseng, Jibin Varghese, Hao Wang, Haoxiang Wang, Heng Wang, Ting-Chun Wang, Fangyin Wei, Xinyue Wei, Jay Zhangjie Wu, Jiashu Xu, Wei Yang, Lin Yen-Chen, Xiaohui Zeng, Yu Zeng, Jing Zhang, Qinsheng Zhang, Yuxuan Zhang, Qingqing Zhao, and Artur Zolkowski. Cosmos world foundation model platform for physical ai, 2025. <https://arxiv.org/abs/2501.03575>.
- [28] NVIDIA Cosmos. cosmos-predict2: General-purpose world foundation models for physical ai. <https://github.com/nvidia-cosmos/cosmos-predict2>, 2025. Accessed: YYYY-MM-DD.
- [29] OpenAI. GPT-5 nano, 2025. <https://openai.com/gpt-5/>. Large language model. Release date: August 7, 2025.
- [30] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3982–3992, 2019.
- [31] Stephen Robertson. Understanding inverse document frequency: on theoretical arguments for idf. *Journal of documentation*, 60(5):503–520, 2004.
- [32] Ola Shorinwa, Zhiting Mei, Justin Lidard, Allen Z Ren, and Anirudha Majumdar. A survey on uncertainty quantification of large language models: Taxonomy, open research challenges, and future directions. *ACM Computing Surveys*, 2025.
- [33] Kihyuk Sohn, Honglak Lee, and Xinchen Yan. Learning structured output representation using deep conditional generative models. *Advances in neural information processing systems*, 28, 2015.
- [34] Charles Spearman. The proof and measurement of association between two things. *The American journal of psychology*, 100(3/4):441–471, 1987.
- [35] Suvrit Sra. A short note on parameter approximation for von mises-fisher distributions: and a fast implementation of i s (x). *Computational Statistics*, 27(1):177–190, 2012.

- [36] Zhiyu Tan, Xiaomeng Yang, Luozheng Qin, and Hao Li. Vidgen-1m: A large-scale dataset for text-to-video generation. *arXiv preprint arXiv:2408.02629*, 2024.
- [37] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004.
- [38] Luhuan Wu and Sinead A Williamson. Posterior uncertainty quantification in neural networks using data augmentation. In *International Conference on Artificial Intelligence and Statistics*, pages 3376–3384. PMLR, 2024.
- [39] Cheng Zhang, Judith Bütepage, Hedvig Kjellström, and Stephan Mandt. Advances in variational inference. *IEEE transactions on pattern analysis and machine intelligence*, 41(8):2008–2026, 2018.
- [40] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018.
- [41] Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*, 2019.

# Appendix

## A Appendix

### A.1 Evaluation Setup

We provide additional details on the evaluation setup.

**Metrics.** We consider the following standard video accuracy metrics: structural similarity index measure (SSIM) [37], peak signal-to-noise ratio (PSNR), learned perceptual image patch similarity (LPIPS) [40], and CLIP cosine similarity score. Note that the SSIM, PSNR, and LPIPS primarily assess visual fidelity while the CLIP score captures more semantic information. We take the negative of the LPIPS score to transform it from an error metric to an accuracy metric. To compute the perceptual metrics, we resize all videos spatially to the same dimensions and subsample the longer videos to ensure that all videos have the same duration. For CLIP, we map both the ground-truth video  $v^{gt}$  and all the generated videos  $v_j^i$  to the visual-semantic space using CLIP. We compute the mean of each metric over all generated videos per task, which represents the assigned value of the metric for that task.

### A.2 Proofs

**Proposition 1** (Uncertainty Decomposition). *Define the total predictive uncertainty in the output video as the differential entropy  $h(V | \ell)$  of the distribution  $f_\theta(V | \ell)$ . Then, this quantity can be decomposed as:*

$$h(V | \ell) = h(V | Z) + h(Z | \ell), \quad (3)$$

where  $h(V | Z)$  represents the epistemic uncertainty in  $v$ , and  $h(Z | \ell)$  the aleatoric uncertainty.

*Proof.* The entropy of a random variable quantifies its associated uncertainty. Given the probability distribution  $f_\theta(V | \ell)$ , we find its entropy by:

$$h(V | \ell) = - \int_{v \in \mathcal{V}} f_\theta(V | \ell) \log(f_\theta(V | \ell)) dv \quad (11)$$

$$= - \int_{v \in \mathcal{V}} \int_{z \in \mathcal{Z}} p(V | z) p(z | \ell) \log(p(V | z) p(z | \ell)) dz dv, \quad (12)$$

where we incorporate the latent state generation step introduced in [Equation \(2\)](#). We can then decompose the log terms into two components:

$$h(V | \ell) = - \int_{v \in \mathcal{V}} \int_{z \in \mathcal{Z}} p(V | z) p(z | \ell) (\log(p(V | z)) + \log(p(z | \ell))) dz dv \quad (13)$$

$$\begin{aligned} &= - \left( \int_{z \in \mathcal{Z}} p(z | \ell) \int_{v \in \mathcal{V}} p(V | z) \log(p(V | z)) dv dz \right) \\ &\quad - \left( \int_{z \in \mathcal{Z}} \left( \int_{v \in \mathcal{V}} p(V | z) dv \right) p(z | \ell) \log(p(z | \ell)) dz \right), \end{aligned} \quad (14)$$

where [Equation \(14\)](#) applies the Fubini-Tonelli theorem. We note that each term of [Equation \(14\)](#) is an entropy itself:

$$h(V | \ell) = - \left( \int_{z \in \mathcal{Z}} p(z | \ell) h(V | Z = z) dz \right) - \left( \int_{z \in \mathcal{Z}} p(z | \ell) \log(p(z | \ell)) dz \right) \quad (15)$$

$$= h(V | Z) + h(Z | \ell). \quad (16)$$

We recognize that the first term  $h(V | Z)$  eliminates uncertainty in prompt ambiguity, and thus signifies the epistemic uncertainty in video generation. On the other hand, the second term  $h(Z | \ell)$  is independent of the video model, but rather only depends on the vagueness of the input prompt, signifying aleatoric uncertainty.  $\square$