

# A Hierarchical Variational Neural Uncertainty Model for Stochastic Video Prediction

Moitreya Chatterjee<sup>1\*</sup> Narendra Ahuja<sup>1</sup> Anoop Cherian<sup>2\*</sup>

<sup>1</sup>University of Illinois at Urbana-Champaign Champaign, IL 61820, USA

<sup>2</sup>Mitsubishi Electric Research Laboratories, Cambridge, MA 02139

metro.smiles@gmail.com n-ahuja@illinois.edu cherian@merl.com

## Abstract

Predicting the future frames of a video is a challenging task, in part due to the underlying stochastic real-world phenomena. Prior approaches to solve this task typically estimate a latent prior characterizing this stochasticity, however do not account for the predictive uncertainty of the (deep learning) model. Such approaches often derive the training signal from the mean-squared error (MSE) between the generated frame and the ground truth, which can lead to sub-optimal training, especially when the predictive uncertainty is high. Towards this end, we introduce Neural Uncertainty Quantifier (NUQ) - a stochastic quantification of the model’s predictive uncertainty, and use it to weigh the MSE loss. We propose a hierarchical, variational framework to derive NUQ in a principled manner using a deep, Bayesian graphical model. Our experiments on four benchmark stochastic video prediction datasets show that our proposed framework trains more effectively compared to the state-of-the-art models (especially when the training sets are small), while demonstrating better video generation quality and diversity against several evaluation metrics.

## 1. Introduction

Extrapolating the present into the future is a task essential to predictive reasoning and planning. When artificial intelligence systems are deployed to work side-by-side with humans, it is critical that they reason about their visual context and generate plausible futures so that they can anticipate the potential needs of humans or catastrophic risks and be better equipped. Such a visual future generation framework could also benefit applications such as video surveillance [59], human action recognition and forecasting [51, 57] as well as simulation of real-world scenarios to train robot learning algorithms, including autonomous driving [28]. However, such applications have a high element of stochasticity, which

\*Equal contribution.

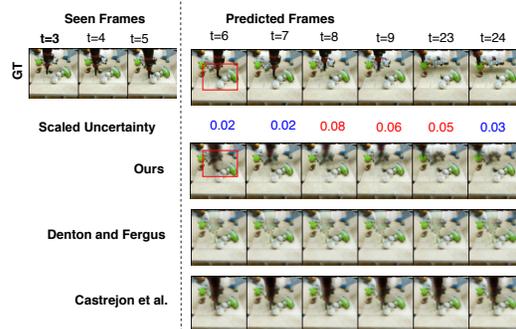


Figure 1. Qualitative results vis-à-vis state-of-the-art video prediction baselines using the proposed NUQ framework on the BAIR Push dataset [15], trained using only 2,000 samples (rather than the full 40K samples). Regions with high motion are shown by a red box. Also shown is an estimate of the per-frame scaled uncertainty estimated by our model. Note that the robotic arm changes direction at  $t = 8$ , which is reflected in the predicted uncertainty.

makes this prediction task challenging.

The resurgence of deep neural networks, especially the advent of generative adversarial networks [20], has enabled significant progress in the development of frameworks for generating visual data, such as images [30]. While, temporally-evolving extensions of such image generation techniques have shown benefits in artificially producing video sequences for deterministic visual contexts [56, 58, 19, 37, 29], they usually fail to model real-world sequences that are often highly stochastic.

Several recent works in video generation, thus design modules to factor in data stochasticity while making predictions [39, 2, 13, 8]. Specifically, such methods assume a latent stochastic prior, from which random samples are drawn, in order to generate future frames. In Babaeizadeh *et al.* [2], this stochastic prior is assumed to follow a fixed normal distribution, which is sampled at every time step, while Denton and Fergus [13], learn this prior from data. The latter’s key insight is to use a variational posterior to guide the learning of the prior to produce the sufficient statistics of the normal distribution governing the prior. Such stochastic

methods typically employ a deterministic decoder (a neural network) that combines an embedding of the visual context and a random sample from the stochastic prior to generate a future video frame. The variance in this prior accounts for the stochasticity underlying the data. To train such models, the mean-squared error (MSE) is then minimized by comparing the predictions against the true video frames.

Nonetheless existing stochastic methods have largely ignored the predictive uncertainty (aleatoric uncertainty) [31] of the models, which might adversarially impact downstream tasks that leverage these predictions. From a machine learning stand point, ignoring the predictive uncertainty might lead to the model being unnecessarily penalized (via the MSE), even if it makes a very uncertain prediction that ends up being different from the ground-truth. This can destabilize the training of the underlying neural networks, leading to slower convergence or requiring larger training data. This is of importance because such data might be expensive or sometimes even difficult to collect (e.g., predicting the next human actions in instruction videos, or a rare traffic incident), and thus effective training with limited data is essential.

In this work, we rise up to these challenges by quantifying the predictive uncertainty of a stochastic frame prediction model and using it to calibrate its training objective. In particular a stochastic estimate of the predictive uncertainty, derived from the latent space of the model, is used to weigh the MSE. That is, when the uncertainty is high, the MSE is down-weighted proportionately, and vice versa; thereby regularizing the backpropagation gradients to train the frame generation module. Moreover, this uncertainty estimate can be used for downstream tasks, such as for example, regulating the maneuvers in autonomous driving [28, 57]. We call our scheme, *Neural Uncertainty Quantifier* (NUQ).

We observe that the weight on the MSE that NUQ introduces, basically amounts to the variance of the normal distribution governing the generated output. Thus, an obvious consideration would be to estimate the variance directly from the output. However, this may be cumbersome due to the very high dimensionality of the output space (order of the number of pixels). We instead, choose to derive it from the variance of the latent space prior, which has far fewer dimensions. Specifically, NUQ leverages a variational, deep, hierarchical, graphical model to bridge the variance of the latent space prior and that of the output. Our framework is trained end-to-end. Sample generations by our framework is shown in Figure 1. In addition, inspired by the recent successes of generative adversarial networks [20, 37, 39], we propose a variant of our framework that uses a novel sequence discriminator, in an adversarial setting. This discriminator module helps to constrain the space of possible output frames, while enforcing motion regularities in the generated videos.

To empirically verify our intuitions, we present experi-

ments on a synthetic (Stochastic Moving MNIST [13]) and three challenging real world datasets: KTH-Action [48], BAIR push [15], and UCF-101 [49] for the task of future frame generation. Our results show that our framework converges faster than prior stochastic video generation methods, and leads to state-of-the-art video generation quality, even when the dataset size is small, while exhibiting generative diversity in the predicted frames.

Below, we summarize the main contributions of this paper:

1. We present *Neural Uncertainty Quantifier* (NUQ), a deep, Bayesian network that learns to estimate the predictive uncertainty of stochastic frame generation models, which can be leveraged to control the training updates, for faster and improved convergence of predictive models.
2. We propose a novel, hierarchical, variational training scheme that allows for incorporating problem-specific knowledge into the predictions via hyperpriors on the uncertainty estimate.
3. Experimental results demonstrate our framework’s better video generation and faster training capabilities, even with small training sets compared to recent state-of-the-art methods on stochastic video generation tasks, across multiple datasets.

## 2. Related Work

Early works in video frame prediction mostly resorted to end-to-end deterministic architectures [44, 50, 17]. Ranzato *et al.* [44] proposed to divide frames into patches and extrapolate their evolution in time. Srivastava *et al.* [50] use image encoders with pre-trained weights to encode the frames. ContextVP [7] and PredNet [41] leverage Convolutional LSTMs [62] for video prediction. Fragkiadaki *et al.* [19] proposes pose extrapolation using LSTMs. More recent approaches [37, 64] seek to predict frames bidirectionally (future and past), during training. However, the inherent deterministic nature of such models [29] often becomes a bottleneck to their performance. Instead, we seek to investigate approaches that allow modeling of the underlying stochasticity in the data while generating an assessment of the model’s predictive uncertainty.

Stochastic approaches constitute a recently emerging and one of the most promising classes of video prediction methods [39, 2, 13]. These approaches model the data stochasticity using a latent prior distribution and are thus readily generalizable to real-world scenarios. Popular among them are STORNs [5], VRNNs [11], SRNNs [18], and DMMs [35]. SV2P [2] is a more recent method that uses a single set of stochastic latent variables that are assumed to follow a fixed prior distribution. Denton and Fergus [13] improve upon

SV2P [2] by allowing the prior distribution to be adapted at every time step by casting the prior as a trainable neural network. Their method is shown to achieve superior empirical performance, thus underlining the importance of learning to model data stochasticity. We also note that generative models have recently been adapted to incorporate stochastic information through a hierarchical latent space [53, 54]. Such networks have also been applied to frame prediction tasks [8]. None of these approaches however, explore the effectiveness of modeling the predictive uncertainty. While, technically it might be possible that the stochastic modules in these prior approaches can learn to quantify this uncertainty implicitly, it may need longer training periods or larger datasets. Instead we show that explicitly incorporating the predictive uncertainty into the learning objective, via a hierarchical, variational framework improves training and inference.

Another line of work in frame prediction seeks to decouple the video into static and moving components [55, 14, 40, 26, 61, 21]. Some of these approaches are deterministic, others stochastic. Denton *et al.* [14] extracts content and pose information for this purpose. Villegas *et al.* [55] adopt a multiscale approach towards frame prediction which works by building a model of object motion, however they require supervisory information, such as annotated pose, during training. Ye *et al.* [63] propose a compositional approach to video prediction by stitching the motion of individual entities. While promising, their approach relies on auxiliary information such as spatial locations of the entities, and as a result, is difficult to generalize. Jin *et al.* [29] investigates decoupling in the frequency space, however they do not model the data stochasticity explicitly. Hsieh *et al.* [27] describes a similar approach by modeling the motion and appearance of each object in the video, but without requiring any auxiliary information. Different from these set of approaches, our proposed framework models frames holistically and is thus agnostic to the video content.

Modeling the predictive uncertainty in deep networks has garnered significant attention lately [6, 12, 36, 42, 52]. Some of these works [1, 36, 42] investigate it in a classification setting, while some others [6, 24] in the context of regression. Uncertainty has also recently been explored in the context of generative models [38, 43, 60]. However, predictive uncertainty modeling in the context of frame prediction has remained largely unexplored. NUQ attempts to fill this gap.

### 3. Background

Suppose  $\mathbf{x}_{1:T} := \langle \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T \rangle$  denotes a sequence of random variables, each  $\mathbf{x}_t$  representing a video frame at time step  $t$ . Assuming we have access to a few initial frames  $\mathbf{x}_{1:F}$ , to set the visual context (where  $1 \leq F < T$ ), our goal is to generate the rest of the frames  $\mathbf{x}_{F+1}$  onwards autoregressively, i.e., conditioned on the seen frames and what has been generated hitherto. This task amounts to

finding a prediction model  $p_\theta(\cdot)$ , parameterized by  $\theta$ , that minimizes the expected negative log-likelihood.

When unknown factors of variation are involved in the data generation process, a deterministic predictive model is insufficient. A standard way to incorporate stochasticity is by assuming the generated frames are in turn conditioned on a latent prior model  $p(\mathbf{z}_t)$ ; i.e.,  $\mathbf{z}_t \sim p(\mathbf{z}_t)$ ,  $\mathbf{x}_t \sim p_\theta(\mathbf{x}_t | \mathbf{x}_{1:t-1}, \mathbf{z}_t)$ . Specifically, the stochasticity in the generative process is characterized by the variance in  $p(\mathbf{z}_t)$ , that produces diversity in  $\mathbf{z}_t \sim p(\mathbf{z}_t)$ . Diversity among predicted frames emerges as a result of this variance.

A well-known problem with the use of such latent stochastic priors is the intractability that it brings into the estimation of the *evidence* or the log-partition function:  $p(\mathbf{x}_t | \mathbf{x}_{1:t-1}) = \int_{\mathbf{z}_t} p_\theta(\mathbf{x}_t | \mathbf{x}_{1:t-1}, \mathbf{z}_t) p(\mathbf{z}_t) d\mathbf{z}_t$ . This problem is typically avoided by casting this estimation in an encoder-decoder setup, where the encoder embeds  $\mathbf{x}_{1:t}$  as  $\mathbf{z}_t \sim p(\mathbf{z}_t | \mathbf{x}_{1:t})$ , while the decoder outputs  $\mathbf{x}_t \sim p_\theta(\mathbf{x}_t | \mathbf{x}_{1:t-1}, \mathbf{z}_t)$ . In order to train efficiently, access to a variational posterior  $q_\phi(\mathbf{z}_t | \mathbf{x}_{1:t})$  – that approximates the true posterior  $p(\mathbf{z}_t | \mathbf{x}_{1:t})$  of the encoder – is assumed. Using this approximate posterior, learning the model parameters  $\theta$  and  $\phi$  amounts to maximizing the variational lower bound,  $\mathcal{L}_{\theta, \phi}$  [33]:

$$\begin{aligned} \log p(\mathbf{x}_t | \mathbf{x}_{1:t-1}) &\geq \mathcal{L}_{\theta, \phi}, \text{ where} \\ \mathcal{L}_{\theta, \phi} &:= \int_{\mathbf{z}_t} q_\phi(\mathbf{z}_t | \mathbf{x}_{1:t}) \log p_\theta(\mathbf{x}_t | \mathbf{x}_{1:t-1}, \mathbf{z}_t) d\mathbf{z}_t \\ &\quad - \int_{\mathbf{z}_t} q_\phi(\mathbf{z}_t | \mathbf{x}_{1:t}) \log \frac{q_\phi(\mathbf{z}_t | \mathbf{x}_{1:t})}{p(\mathbf{z}_t)} d\mathbf{z}_t \end{aligned} \quad (1)$$

From the definition, this amounts to:

$$\begin{aligned} \mathcal{L}_{\theta, \phi} &= \mathbb{E}_{q_\phi(\mathbf{z}_t | \mathbf{x}_{1:t})} \log p_\theta(\mathbf{x}_t | \mathbf{x}_{1:t-1}, \mathbf{z}_t) - \\ &\quad \text{KL}(q_\phi(\mathbf{z}_t | \mathbf{x}_{1:t}) || p(\mathbf{z}_t)), \text{ for } t > F. \end{aligned} \quad (2)$$

Leveraging the re-parametrization trick ([33]) allows efficient optimization of the likelihood loss in Eq. 2, permitting us to learn the parameters  $\theta$  and  $\phi$ . Note that the expectation term in Eq. 2 boils down to a standard MSE over all predicted frames  $\mathbf{x}_{F+1:T}$  in the training set when the  $p_\theta(\cdot)$  term is assumed to follow a Gaussian distribution with an isotropic constant variance. In this setting, the KL divergence in Eq. 2 acts as a regularizer on  $q_\phi(\cdot)$  so that this posterior does not just copy an encoding of  $\mathbf{x}_t$  available to it as  $\mathbf{z}_t$ , instead captures the density of a latent distribution that is useful to the prediction model in maximizing the first term in Eq. 2.

In conditional variational autoencoders [2, 33], the latent prior  $p(\mathbf{z}_t)$  is typically assumed to be  $\mathcal{N}(0, 1)$  - a choice that can be sub-optimal. A better approach is perhaps to learn this prior so that the stochasticity of the future frame can be guided by the data itself. To this end, Denton and Fergus [13] suggests a learned stochastic prior model  $p(\mathbf{z}_t) = p_\psi(\mathbf{z}_t) := p_\psi(\mathbf{z}_t | \mathbf{x}_{1:t-1})$ , parametrized by  $\psi$ , which is learned by minimizing its divergence from the

variational posterior  $q_\phi(\cdot)$  through the KL-term in Eq. 2. As the posterior  $q_\phi(\cdot)$  has access to the current input sample  $\mathbf{x}_t$ , it can guide the prior (which does not have access to  $\mathbf{x}_t$ , but only to  $\mathbf{x}_{1:t-1}$ ) to produce a distribution on  $\mathbf{z}_t$  that mimics the posterior (and hence we can discard the posterior at test time). Thus, the training-time sampling pipeline is given by:  $\mathbf{z}_t \sim q_\phi(\mathbf{z}_t|\mathbf{x}_{1:t})$ ,  $\mathbf{x}_t \sim p_\theta(\mathbf{x}_t|\mathbf{x}_{1:t-1}, \mathbf{z}_t)$ , and  $p_\psi \stackrel{d}{\leftarrow} q_\phi$  (matching in distribution).

While learning the stochastic prior  $p_\psi(\cdot)$  allows for characterizing the data stochasticity, the model’s predictive uncertainty remains unaccounted for. Our hierarchical framework for estimating this uncertainty, follows a two-step process. The first is the estimation of the learned prior,  $p_\psi(\cdot)$ . The key idea in the second step is to leverage the variance in the learned prior  $p_\psi(\cdot)$  to estimate this uncertainty. Since the prior estimation network,  $p_\psi(\cdot)$ , is retained both during training and inference (unlike the posterior), this permits its usage for downstream tasks, during inference.

## 4. Proposed Method

As alluded to above, we seek to control the prediction model using the uncertainty estimated directly from the stochastic prior. Subsequently, we assume the prediction model consists of an LSTM,  $f_\theta$ , with weights  $\theta$  such that:

$$\hat{\mathbf{x}}_t = f_\theta(\mathbf{x}_{1:t-1}, \mathbf{z}_t) := f_\theta(\mathbf{x}_{t-1}, \mathbf{z}_t; \mathbf{h}_{t-1}^\theta), \quad (3)$$

where  $\hat{\mathbf{x}}_t$  denotes the  $t^{\text{th}}$  generated frame and  $\mathbf{h}_{t-1}^\theta$  captures the internal states of the LSTM. The generated frame  $\hat{\mathbf{x}}_t$  is then sent through the likelihood model to compute the MSE. With this setup, we are now ready to introduce our neural uncertainty quantifier (NUQ). Figure 2 provides an overview of our framework.

### 4.1. Neural Uncertainty Quantifier

As is standard practice, let us assume the data likelihood model  $p_\theta(\mathbf{x}_t|\mathbf{x}_{1:t-1}, \mathbf{z}_t) \sim \mathcal{N}(\mathbf{x}_t, \frac{1}{b_t})$ , where  $b_t$  denotes the precision (inverse variance) of our prediction model, where  $b_t > 0$ ,  $b_t \in \mathbb{R}$ . Denton and Fergus [13] assumes  $b_t$  to be an isotropic constant, such that the negative log-likelihood of the predicted frame  $\hat{\mathbf{x}}_t$  boils down to computing the  $\ell_2$ -loss. This reduces Eq. 2 to become the evidence lower bound (ELBO) [33]:

$$\mathcal{L}_{\theta, \phi, \psi} := \sum_{t=F+1}^T \frac{1}{2} \|\hat{\mathbf{x}}_t - \mathbf{x}_t\|^2 + \text{KL}(q_\phi(\mathbf{z}_t|\mathbf{x}_{1:t}) \| p_\psi(\mathbf{z}_t|\mathbf{x}_{1:t-1})). \quad (4)$$

Our key insight to the proposed uncertainty measure arises from the observation that the  $\ell_2$ -norm term in Eq. 4 does not include any dependency on the uncertainty associated with the prediction of  $\hat{\mathbf{x}}_t$ . Note that there are two

extreme situations when this loss is large that impacts effective training: (i) when there is no uncertainty associated with the generation of  $\hat{\mathbf{x}}_t$ , however the prediction model is not trained well, such that  $\hat{\mathbf{x}}_t$  does not match  $\mathbf{x}_t$ , and (ii) when there is uncertainty involved in the generative model such that the generated  $\hat{\mathbf{x}}_t$ , while plausible given the context, is different from  $\mathbf{x}_t$ . Thus, the key research question for effective model training becomes: how can we equip the prediction model to differentiate these situations? Our solution is to directly condition the prediction model with the uncertainty derived from the prior  $p_\psi(\mathbf{z}_t|\mathbf{x}_{1:t-1})$ , so that when the stochasticity is high for the generated frames, the  $\ell_2$ -loss term is weighed down such that the gradients computed on this term will have a lesser impact in updating the weights of the neural network; thereby stabilizing the training.

Suppose our prior  $p_\psi(\mathbf{z}_t|\mathbf{x}_{1:t-1})$  is a normal distribution  $\mathcal{N}(\boldsymbol{\mu}_t^z, \boldsymbol{\Sigma}_t^z)$ , with parameters  $\boldsymbol{\mu}_t^z$ , the mean, and  $\boldsymbol{\Sigma}_t^z$  the covariance matrix - capturing the predictive uncertainty, in the latent space. For better characterization of this prior model, we assume it to be implemented as an LSTM  $g_\psi$  with weights  $\psi$  such that  $(\boldsymbol{\mu}_t^z, \boldsymbol{\Sigma}_t^z) = g_\psi(\mathbf{x}_{t-1}; \mathbf{h}_{t-1}^\psi)$ , where  $\mathbf{h}_{t-1}^\psi$  denotes the hidden state. Similarly, we assume the posterior  $q_\phi(\mathbf{z}_t|\mathbf{x}_{1:t})$  is normally-distributed:  $\mathcal{N}(\boldsymbol{\mu}_t^q, \boldsymbol{\Sigma}_t^q)$ , and is implemented using an LSTM  $l_\phi(\mathbf{x}_t; \mathbf{h}_t^\phi)$  with weights  $\phi$  and hidden state  $\mathbf{h}_t^\phi$ . This leads us to the following sampling pipeline:

$$\hat{\mathbf{x}}_t = f_\theta(\mathbf{x}_{t-1}, \mathbf{z}_t; \mathbf{h}_{t-1}^\theta) \sim \mathcal{N}(\mathbf{x}_t, \frac{1}{b_t}), \quad (5)$$

$$\mathbf{z}_t|\mathbf{x}_{1:t-1} \sim \mathcal{N}(\boldsymbol{\mu}_t^z, \boldsymbol{\Sigma}_t^z); (\boldsymbol{\mu}_t^z, \boldsymbol{\Sigma}_t^z) = g_\psi(\mathbf{x}_{t-1}; \mathbf{h}_{t-1}^\psi), \quad (6)$$

$$\mathbf{z}_t|\mathbf{x}_{1:t} \sim \mathcal{N}(\boldsymbol{\mu}_t^q, \boldsymbol{\Sigma}_t^q); (\boldsymbol{\mu}_t^q, \boldsymbol{\Sigma}_t^q) = l_\phi(\mathbf{x}_t; \mathbf{h}_t^\phi), \quad (7)$$

Using this setup, we are now ready to present our hierarchical, generative, variational model for uncertainty estimation, an overview of which is shown in Figure 2. To set the stage, let us assume the precision is sampled from the distribution  $p(b_t|\mathbf{x}_{1:t-1})$ . Then, we can rewrite the log-likelihood in Eq. 1 by including the precision distribution as:

$$\int_{b_t, z_t} \log p(\mathbf{x}_t|b_t, \mathbf{z}_t, \mathbf{x}_{1:t-1}) + \log p(\mathbf{z}_t|\mathbf{x}_{1:t-1}) + \log p(b_t|\mathbf{x}_{1:t-1}) db_t dz_t \quad (8)$$

The above integral is intractable. Hence, we approximate it by sampling  $b_t$  and  $\mathbf{z}_t$ . Note that the first two terms taken together is essentially the left-hand side of Eq. 1, except with the additional conditioning on  $b_t$ . Using the variational lower bound [33], like in Eq. 1, we have :

$$\log p(\mathbf{x}_t|b_t, \mathbf{x}_{1:t-1}) \geq \mathbb{E}_{q_\phi(\mathbf{z}_t|\mathbf{x}_{1:t})} \log p_\theta(\mathbf{x}_t|\mathbf{x}_{1:t-1}, \mathbf{z}_t, b_t) - \text{KL}(q_\phi(\mathbf{z}_t|\mathbf{x}_{1:t}) \| p_\psi(\mathbf{z}_t|\mathbf{x}_{1:t-1})), \text{ for } t > F. \quad (9)$$

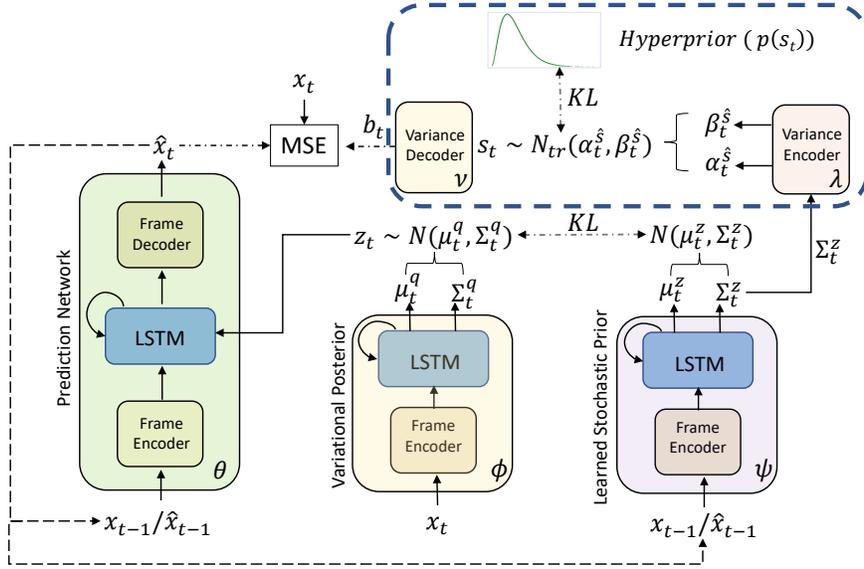


Figure 2. Overview of our approach.

Please see the supplementary for the derivation.

As stated before, we seek to connect the uncertainty  $\Sigma_t^z$  in the latent prior  $p_\psi(z_t|x_{1:t-1})$  to the precision  $b_t$ . We accomplish this via a variational encoder-decoder network. Such a formulation permits the flexibility of introducing customized prior distributions on its latent space. During training, the encoder component of this network,  $\zeta_\lambda(\cdot)$ , with parameters  $\lambda$ , takes as input  $\Sigma_t^z$ , and produces the sufficient statistics of the posterior distribution  $q_\lambda(\cdot)$  governing the latent space of this network, while the decoder  $\tau_\nu(\cdot)$ , with parameters  $\nu$  draws a sample  $s_t$ , from this distribution, and decodes it to generate  $b_t$ , with a distribution on  $b_t$  denoted by  $p_\nu(\cdot)$ . This sampling scheme is described as follows:

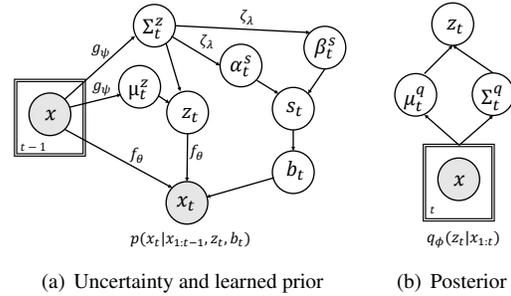
$$\begin{aligned} b_t &\sim p_\nu(b_t|s_t, \mathbf{x}_{1:t-1}) = p_\nu(b_t|s_t, \Sigma_t^z), \\ s_t &\sim q_\lambda(s_t|\mathbf{x}_{1:t-1}) = q_\lambda(s_t|\Sigma_t^z) \end{aligned} \quad (10)$$

In order to provide appropriate regularization for the latent space distribution,  $q_\lambda(\cdot)$ , we assume a manually-defined hyper-prior distribution governing the latent space of this module denoted  $p(s_t)$ . Let the hyper prior  $p(s_t) \sim D_\gamma(\alpha_t^s, \beta_t^s)$ , with parameters  $\alpha_t^s, \beta_t^s$  chosen by the user and let  $q_\lambda(s_t|\Sigma_t^z) \sim D_\lambda(\alpha_t^{\hat{s}}, \beta_t^{\hat{s}})$ , where the parameters  $\alpha_t^{\hat{s}}, \beta_t^{\hat{s}}$  are estimated by the encoder network  $\zeta_\lambda(\cdot)$ . With this setup, analogous to Eq. 2, we obtain the following variational lower bound on the likelihood of  $b_t$ :

$$\begin{aligned} \log p(b_t|\mathbf{x}_{1:t-1}) &\geq \mathbb{E}_{q_\lambda(s_t|\Sigma_t^z)} \log p_\nu(b_t|s_t, \Sigma_t^z) - \\ &\text{KL}(q_\lambda(s_t|\Sigma_t^z) \| p(s_t)), \text{ for } t > F \end{aligned} \quad (11)$$

Please see the supplementary for the derivation. Plugging Eq. 9 and Eq. 11 in Eq. 8, we have the following:

$$\begin{aligned} &\mathbb{E}_{q_\phi(z_t|\mathbf{x}_{1:t})} \log p_\theta(\mathbf{x}_t|\mathbf{x}_{1:t-1}, \mathbf{z}_t, b_t) - \\ &\quad \text{KL}(q_\phi(z_t|\mathbf{x}_{1:t}) \| p_\psi(z_t|\mathbf{x}_{1:t-1})) + \\ &\mathbb{E}_{q_\lambda(s_t|\mathbf{x}_{1:t-1})} \log p_\nu(b_t|s_t, \Sigma_t^z) - \text{KL}(q_\lambda(s_t|\Sigma_t^z) \| p(s_t)) \end{aligned} \quad (12)$$



(a) Uncertainty and learned prior

(b) Posterior

Figure 3. Plate diagrams depicting the graphical model of our NUQ-framework. (a) shows the sampling dependencies between the learned prior and the uncertainty prediction modules, while (b) shows our posterior sampling framework. The plates denote, for example,  $t-1$  repetitions of the random variable  $\mathbf{x}$  in (a).

Assuming that  $p_\theta(\mathbf{x}_t|\mathbf{x}_{1:t-1}, \mathbf{z}_t, b_t)$  follows a Gaussian distribution  $\mathcal{N}(\mathbf{x}_t, \frac{1}{b_t})$ , along the lines of Eq. 4, leads us to our final ELBO objective, which we minimize:

$$\begin{aligned} \mathcal{L}_{\theta, \phi, \psi, \lambda}^P &= \sum_{t=F+1}^T \frac{1}{2} [b_t \|\hat{\mathbf{x}}_t - \mathbf{x}_t\|^2 - \log b_t] - \\ &\mathbb{E}_{q_\lambda(s_t|\Sigma_t^z)} \log p_\nu(b_t|s_t, \Sigma_t^z) + \\ &\eta_1 \text{KL}(q_\phi(\mathbf{z}_t|\mathbf{x}_{1:t}) \| p_\psi(\mathbf{z}_t|\mathbf{x}_{1:t-1})) + \eta_2 \text{KL}(q_\lambda(s_t|\Sigma_t^z) \| p(s_t)) \end{aligned} \quad (13)$$

where  $\eta_1, \eta_2 \geq 0$  are regularization constants (as suggested in Higgins *et al.* [25]).

Given our setup, a natural choice in a hierarchical Bayesian conjugate sense is to assume  $p(s_t) \sim \Gamma(\alpha_t^s, \beta_t^s)$ , the gamma distribution, which is a conjugate prior for the precision. Unfortunately, however using the gamma distribution for the posterior does not permit the reparameterization trick [34, 33], which is essential for making sampling-based networks differentiable. While one may resort to approximations to the gamma prior such as using implicit gradients [16] or generalized re-parameterization techniques [47], these ap-

proaches can be computationally expensive. Instead, we propose to approximate it by a truncated normal distribution  $\mathcal{N}_{tr}(\alpha_t^{\hat{s}}, \beta_t^{\hat{s}})$  (which is amenable to re-parametrization), where now  $\alpha_t^{\hat{s}} (\geq 0)$  and  $\beta_t^{\hat{s}}$  correspond to the mean and the standard deviation of the truncated normal, respectively and are estimated by the encoder network  $\zeta_\lambda(\cdot)$ . In practice, the truncation is effected through rejection sampling [45]; i.e., we sample from a normal distribution, and reject samples if they are negative. Empirically, we find this choice of the hyper prior (being a gamma distribution) and our truncated-normal posterior combination to be beneficial. Thus, the KL divergence in Eq. 11 will eventually promote the network  $\zeta_\lambda(\cdot)$  to produce the sufficient statistics of a truncated-normal distribution which will closely approximate the true gamma hyper-prior governing  $p(s_t)$ . Additionally, since  $s_t$  is a scalar ( $\in \mathbb{R}$ ), thus  $b_t (= \frac{1}{s_t})$  is directly sampled from  $q_\lambda(\cdot)$  rather than through the decoder network,  $\tau_\nu(\cdot)$ . Figure 3 presents a plate diagram of our proposed hierarchical graphical model.

## 4.2. Sequence Discriminator

Inspired by the success of generative adversarial networks in generating realistic images and realistic object motions [20, 20, 22, 4, 39, 37], as well as the synergy that GANs bring about in improving the quality of other generative models [10, 23, 3], we propose to integrate NUQ with a sequence discriminator, where the generated frame sequences are input to a discriminator that checks for their realism and motion coherence. Different from prior approaches that employ GANs for future frame prediction [39, 37], our discriminator ( $D_w$ ) is a recurrent neural network with weights  $w$ . It takes as input  $k$  contiguous frames with image dimensions  $\delta_h \times \delta_w$ , and produces a non-negative score, denoting the probability of that sequence being real or fake. Thus,  $D_w : \mathbb{R}^{\delta_h \times \delta_w \times k} \rightarrow \mathbb{R}_+$ . Suppose,  $\mathbf{y}_{1:k} \subset \mathbf{x}_{1:T}$  represents  $k$  contiguous frames starting at a random time step from video sequences  $\mathbf{x}_{1:T}$  in the training dataset  $\mathcal{X}$ . If  $\hat{\mathbf{x}}_{t-k+1:t}$  represents a sequence of  $k$  generated frames, then our discriminator loss is given by:

$$\mathcal{L}_w^D := - \sum_{t=F+1}^T \mathbb{E}_{\mathbf{y}_{1:k} \sim \mathcal{X}} \log [D_w(\mathbf{y}_{1:k})] + \quad (14)$$

$$\mathbb{E}_{\hat{\mathbf{x}}_{t-k+1:t} \sim p(\hat{\mathbf{x}}_{t-k+1:t} | \mathbf{x}_{1:t-k}, \mathbf{z}_{1:t})} \log [1 - D_w(\hat{\mathbf{x}}_{t-k+1:t})],$$

where, the discriminator is trained to distinguish between the generated (with label zero) and real (with label one) input sequences, by minimizing  $\mathcal{L}_w^D$ , while the generator tries to maximize it. Combining the ELBO loss in Eq. 13 with the generator loss, we have our modified training loss for this variant, given by  $\mathcal{L} = \mathcal{L}_{\theta, \phi, \psi, \lambda}^P - \gamma \mathcal{L}_w^D$ , where  $\gamma > 0$  is a small regularization constant. For both variants (Eq. 13 or Eq. 14), we optimize the final objective using ADAM [32].

## 5. Experiments

In this section, we empirically validate the efficacy of NUQ on challenging real-world and synthetic datasets.

**Datasets:** We conduct experiments on four standard stochastic video prediction datasets, namely (i) Stochastic Moving MNIST (SMMNIST) [13, 10], (ii) BAIR Robot Push [15], (iii) KTH-Action [48], and (iv) UCF-101 [49]. In SMMNIST, a hand-written digit moves in rectilinear paths within a  $48 \times 48$  pixel box, bouncing off its walls, where the post bounce movement directions are stochastic. The dataset has a test set size of 1,000 videos [10]. The BAIR Push Dataset [15] consists of  $64 \times 64$  pixel videos featuring highly stochastic motions of a Sawyer robotic arm pushing objects on a table. This dataset has 257 test samples using the split of Denton and Fergus [13]. The KTH Action Dataset [48] is a small dataset of  $64 \times 64$  pixel videos containing a human performing various actions (walking, jogging, etc.), captured in a controlled setting with a static camera. The test set for this dataset consists of 476 videos. Finally, the UCF-101 Dataset [49] is a dataset of videos, resized to  $64 \times 64$  containing 101 common human action categories (such as pushups, cricket shot, etc.), spanning both indoor and outdoor activities. The test set for this dataset consists of 1,895 videos. We hypothesize that by incorporating the predictive uncertainty, NUQ undergoes more efficient training updates and can thus train with fewer samples, efficiently. We therefore conduct experiments with varying training set sizes for some of these datasets.

**Baselines and Evaluation Metrics:** To carefully evaluate the performance improvement brought about by incorporating our uncertainty estimation method into a stochastic video generation framework, we choose three competitive and closely-related state-of-the-art methods within the stochastic video prediction realm as baselines, namely: (i) Denton and Fergus [13], (ii) Castrejon *et al.* [8], and (ii) Hsieh *et al.* [27]. At test time, we follow the standard protocol of generating 100 sequences for all models and report performances on sequences that matches best with the ground truth [13]. To quantify the generation quality, we use standard evaluation metrics: (i) per-frame Structural Similarity (SSIM) ([9]), (ii) Peak Signal to Noise Ratio (PSNR), and (iii) Learned Perceptual Image Patch Similarity (LPIPS) [65] - with a VGG backbone. We report the average scores on these metrics across all predicted frames.

**Experimental Setup:** For SMMNIST, BAIR Push, and KTH Action, we train all methods with 5 seen and 15 unseen frames, while at test time 20 frames are predicted after the first 5 seen ones. When training with UCF-101, 15 seen and 10 unseen frames are used, while at test time the number of unseen frames is set to 15. For the baseline methods [13, 8, 27], we use the publicly available implementations from the authors. To ensure our proposed NUQ-framework is similar in learning capacity, we use the

Table 1. SSIM, PSNR, and LPIPS scores on the test set for different datasets after @1, @5, and @Convergence (C) (upto 150 epochs) epochs of training with varying training set sizes. [Key: Best results in **bold** and second-best in **blue**.]

Dataset: SMMNIST	SSIM ↑			PSNR ↑			LPIPS ↓		
	@1	@5	@C	@1	@5	@C	@1	@5	@C
<i>Number of training samples - 2,000</i>									
Ours	<b>0.8686</b>	<b>0.8638</b>	<b>0.8948</b>	17.76	<b>18.13</b>	18.14	<b>0.3087</b>	<b>0.2836</b>	<b>0.1803</b>
Ours (w/o discriminator)	0.8599	0.8825	0.8929	<b>17.82</b>	18.07	<b>18.48</b>	0.3283	0.3158	0.1967
Denton and Fergus [13]	0.8145	0.8650	0.8696	17.07	18.05	18.13	0.3429	0.3345	0.2321
Castrejon <i>et al.</i> [8]	0.8564	0.8748	0.8868	17.36	17.98	18.12	0.3392	0.3432	0.2262
Hsieh <i>et al.</i> [27]	0.4538	0.8419	0.8569	11.27	16.40	16.70	0.4370	0.3696	0.2842
<i>Number of training samples - 8,000 (full training set)</i>									
	@1	@5	@C	@1	@5	@C	@1	@5	@C
Ours	<b>0.8524</b>	<b>0.8610</b>	<b>0.9088</b>	<b>17.93</b>	18.14	<b>19.07</b>	<b>0.3787</b>	<b>0.3013</b>	0.1149
Denton and Fergus [13]	0.8154	0.8607	0.8819	17.49	<b>18.22</b>	18.30	0.4061	0.3626	0.2813
Castrejon <i>et al.</i> [8]	0.8640	0.8708	0.8868	17.23	18.06	18.27	0.3939	0.3316	<b>0.1040</b>
Hsieh <i>et al.</i> [27]	0.5328	0.8374	0.8801	11.46	16.65	16.70	0.4217	0.4039	0.2747
<i>Number of training samples - 2,000</i>									
Dataset: BAIR Push	SSIM ↑			PSNR ↑			LPIPS ↓		
	@1	@5	@C	@1	@5	@C	@1	@5	@C
Ours	<b>0.7709</b>	<b>0.8230</b>	<b>0.8314</b>	<b>18.40</b>	<b>19.15</b>	<b>19.26</b>	<b>0.3394</b>	<b>0.2014</b>	<b>0.1574</b>
Denton and Fergus [13]	0.7351	0.7853	0.8196	17.32	17.44	18.49	0.3531	0.3197	0.1725
Castrejon <i>et al.</i> [8]	0.7094	0.7961	0.8221	17.19	17.92	18.79	0.3433	0.2560	0.1742
Hsieh <i>et al.</i> [27]	0.4979	0.7901	0.7989	11.32	15.28	16.00	0.4159	0.3899	0.1891
<i>Number of training samples - 43,264 (full training set)</i>									
	@1	@5	@C	@1	@5	@C	@1	@5	@C
Ours	<b>0.8135</b>	<b>0.8336</b>	<b>0.8460</b>	<b>19.03</b>	<b>19.14</b>	<b>19.31</b>	<b>0.1656</b>	<b>0.1470</b>	0.1296
Denton and Fergus [13]	0.7782	0.8198	0.8328	18.30	18.38	18.81	0.2119	0.1843	0.1499
Castrejon <i>et al.</i> [8]	0.7816	0.8309	0.8437	18.29	18.56	<b>19.59</b>	0.1878	0.1720	<b>0.1181</b>
Hsieh <i>et al.</i> [27]	0.7507	0.8123	0.8323	16.52	16.61	16.61	0.2140	0.1829	0.1713
<i>Number of training samples - 1,911 (full training set)</i>									
Dataset: KTH Action	SSIM ↑			PSNR ↑			LPIPS ↓		
	@1	@5	@C	@1	@5	@C	@1	@5	@C
Ours	<b>0.7990</b>	<b>0.8192</b>	0.8448	<b>22.62</b>	22.89	<b>24.02</b>	<b>0.4309</b>	<b>0.3390</b>	<b>0.2238</b>
Denton and Fergus [13]	0.7028	0.8056	0.8374	21.29	<b>22.93</b>	24.73	0.4621	0.3580	0.2497
Castrejon <i>et al.</i> [8]	0.6345	0.8054	<b>0.8510</b>	21.31	21.12	24.82	0.4513	0.3471	0.2395
Hsieh <i>et al.</i> [27]	0.4647	0.5335	0.7057	11.25	12.32	16.44	0.5189	0.3939	0.2771
<i>Number of training samples - 11,425 (full training set)</i>									
Dataset: UCF-101	SSIM ↑			PSNR ↑			LPIPS ↓		
	@1	@5	@C	@1	@5	@C	@1	@5	@C
Ours	<b>0.7359</b>	<b>0.7636</b>	0.7729	<b>21.25</b>	<b>21.98</b>	<b>22.73</b>	<b>0.3914</b>	<b>0.2865</b>	<b>0.0836</b>
Denton and Fergus [13]	0.6253	0.7540	0.7603	19.35	20.60	21.64	0.3507	0.3006	0.1259
Castrejon <i>et al.</i> [8]	0.6712	0.7555	<b>0.7756</b>	20.58	20.58	22.53	0.3414	0.2965	0.1036
Hsieh <i>et al.</i> [27]	0.6199	0.6800	0.7103	16.65	17.18	18.41	0.3989	0.3239	0.1771

Table 2. SSIM, PSNR, and LPIPS scores on the SMMNIST test set after @1, @5, and @Convergence (C) (upto 150 epochs) epochs of training with alternative formulations of our model using 2,000 training samples. [Key: Best results in **bold**.]

Dataset: SMMNIST	SSIM ↑			PSNR ↑			LPIPS ↓		
	@1	@5	@C	@1	@5	@C	@1	@5	@C
$p(s_t) \sim \text{Uniform}[0, 1]$	0.8173	0.8374	0.8523	17.6	17.95	18.06	0.3442	0.3038	0.198
Estimate $b_t$ from the decoder $p_\theta(\cdot)$	0.7627	0.7628	0.7828	17.54	17.55	17.55	0.3463	0.3259	0.2225
Estimate $b_t$ w/o variance encoder-decoder	0.7450	0.7454	0.7648	16.22	16.53	16.78	0.3469	0.3263	0.2328
NUQ (Ours)	<b>0.8686</b>	<b>0.8638</b>	<b>0.8948</b>	<b>17.76</b>	<b>18.13</b>	<b>18.14</b>	<b>0.3087</b>	<b>0.2836</b>	<b>0.1803</b>

“DCGAN encoder-decoder” architecture for the frames, as in Denton and Fergus [13], for all datasets. Our variance encoder network  $\zeta_\lambda(\cdot)$  is a multi-layer perceptron with 2

layers, and our sequence discriminator is an LSTM with a single 256-d hidden layer. More architectural details are in the Supplementary Materials. We set  $k = 3$  in Eq. 14,

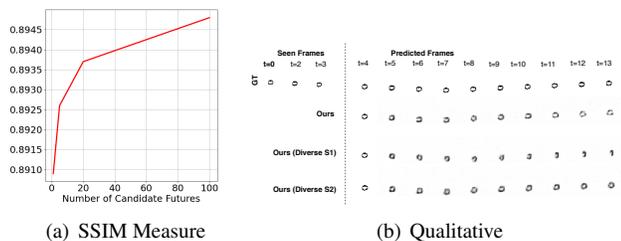


Figure 4. Diversity Results: SSIM score with increasing number of generated futures, per time step, on the SMMNIST test set, as a quantification of output diversity (left) and qualitative generation results (right) using NUQ trained with 2000 training samples.

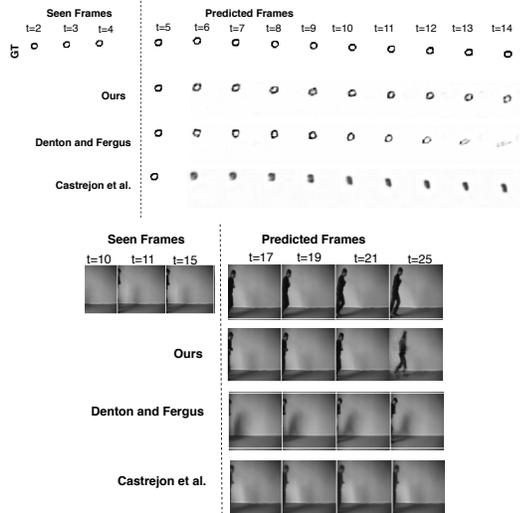


Figure 5. Qualitative results of NUQ against competing baselines on SMMNIST (top) and on KTH-Action (bottom).

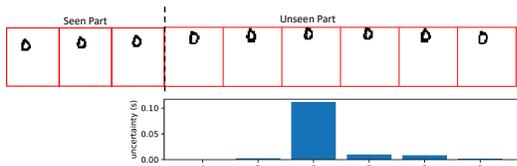


Figure 6. Evolution of scaled uncertainty on the SMMNIST dataset (trained with 2,000 training samples) against time-steps.

Table 3. Human preference scores for samples generated using NUQ versus competing baselines across different datasets,

Dataset (# Training Samples)	Prefer: Ours/ [8]/ [13]
SMMNIST (2000 samples)	89% / 11% / 0 %
BAIR Push (2000 samples)	78% / 22% / 0 %
KTH-Action (1911 samples)	78% / 11% / 11 %

and  $\gamma$  to be 0.00001. We use  $\eta_1 = 0.0001$  and  $\eta_2 = 0.001$  in Eq. 13. Learning rate is set to 0.002 and no scheduling is used. The hyper-parameters for the baseline methods are chosen using a small validation set ( $\sim 5\%$  of training data).

**State-of-the-Art Comparisons:** In Table 1, we report quantitative evaluations of our model versus competing baselines across the four datasets. We observe that both variants of the NUQ outperform recent competitive baselines by wide margins on all measures (upto 10% in SSIM), with the one with

the discriminator being slightly better - underscoring the benefits of adversarial training. While noticeable gains are obtained across the board, NUQ really shines under limited training set sizes. We surmise that this gain is attributable to the failure of the baseline methods in incorporating predictive uncertainty explicitly into the learning objective.

From the table, we also see that NUQ converges faster than other methods both for small and large training set sizes. Figures 1 and 5 show sample generation results from SMMNIST, BAIR Push, and KTH-Action datasets versus competing baselines, trained with 2,000 and 1,911 samples, respectively. From these figures, we see that compared to baseline methods, frames generated by our method are superior at capturing both the appearance and the motion of the object (i.e. digit/robot arm/human) even under limited training data. Human annotators, when presented with a few random sample generations by NUQ versus competing methods, overwhelmingly choose NUQ samples for their realism, as shown in Table 3. Figure 6 shows the evolution of a scaled uncertainty estimate derived from  $\frac{1}{b_t}$  over different frames. The plot shows the increase in uncertainty co-occurs with the bounce of the digit against the boundary, suggesting that the uncertainty is well grounded in the data.

**Alternative Formulations:** Next, we discuss the results of some alternative formulations of our model. We consider: (i) replacing the gamma hyperprior on  $p(s_t)$  with Uniform(0, 1) distribution, (ii) estimating  $b_t$  from the frame decoder  $p_\theta(\cdot)$  by producing a diagonal covariance matrix, and (iii) assuming a deterministic mapping from  $\Sigma_t^z$  to  $b_t$  through a multi-layer perceptron. Table 2 presents the performance of these alternatives on the SMMNIST dataset. From the first row of the table, we see that choosing the Uniform(0, 1) as priors results in suboptimal variants of NUQ. Further, the results also show that either estimating  $b_t$  directly from the decoder  $p_\theta(\cdot)$ , or computing it deterministically from  $\Sigma_t^z$  performs poorly, suggesting that such estimation techniques are not ideal.

**Diversity:** In Figure 4 (a), we plot the average SSIM of NUQ for a set of futures, with increasing cardinality of this set. Our plot shows that the SSIM increases with larger number of futures, suggesting that the possibility of matching with a ground truth future increases with more futures, implying better diversity of our model. Figure 4(b) presents diverse generation results on the SMMNIST dataset by NUQ.

## 6. Conclusions

Recent approaches have demonstrated the need for modeling data uncertainty in video prediction models. However, in this paper we show that the state of the art in such stochastic methods do not leverage the model’s predictive uncertainty to the fullest extent. Indeed, we show that explicitly incorporating this uncertainty into the learning objective via our proposed Neural Uncertainty Quantifier (NUQ) framework,

can lead to faster and more effective model training even with fewer training samples, as validated by our experiments.

## References

- [1] Lynton Ardizzone, Radek Mackowiak, Carsten Rother, and Ullrich Köthe. Training normalizing flows with the information bottleneck for competitive generative classification. *Advances in Neural Information Processing Systems*, 33, 2020. [3](#)
- [2] Mohammad Babaeizadeh, Chelsea Finn, Dumitru Erhan, Roy H Campbell, and Sergey Levine. Stochastic variational video prediction. In *International Conference on Learning Representations*, 2018. [1](#), [2](#), [3](#)
- [3] Jianmin Bao, Dong Chen, Fang Wen, Houqiang Li, and Gang Hua. Cvae-gan: fine-grained image generation through asymmetric training. In *Proceedings of the IEEE international conference on computer vision*, pages 2745–2754, 2017. [6](#)
- [4] Emad Barsoum, John Kender, and Zicheng Liu. Hp-gan: Probabilistic 3d human motion prediction via gan. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 1418–1427, 2018. [6](#)
- [5] Justin Bayer and Christian Osendorfer. Learning stochastic recurrent networks. *arXiv preprint arXiv:1411.7610*, 2014. [2](#)
- [6] Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural network. In *International Conference on Machine Learning*, pages 1613–1622, 2015. [3](#)
- [7] Wonmin Byeon, Qin Wang, Rupesh Kumar Srivastava, and Petros Koumoutsakos. Contextvp: Fully context-aware video prediction. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 753–769, 2018. [2](#)
- [8] Lluís Castrejon, Nicolas Ballas, and Aaron Courville. Improved conditional vrns for video prediction. In *Proceedings of the IEEE International Conference on Computer Vision*, 2019. [1](#), [3](#), [6](#), [7](#), [8](#)
- [9] Sumohana S Channappayya, Alan Conrad Bovik, and Robert W Heath Jr. Rate bounds on ssim index of quantized images. *IEEE Transactions on Image Processing*, 17(9):1624–1639, 2008. [6](#)
- [10] Anoop Cherian, Moitreyee Chatterjee, and Narendra Ahuja. Sound2sight: Generating visual dynamics from sound and context. In *European Conference on Computer Vision (ECCV)*, 2020. [6](#)
- [11] Junyoung Chung, Kyle Kastner, Laurent Dinh, Kratarth Goel, Aaron C Courville, and Yoshua Bengio. A recurrent latent variable model for sequential data. In *Advances in neural information processing systems*, pages 2980–2988, 2015. [2](#)
- [12] William R Clements, Benoît-Marie Robaglia, Bastien Van Delft, Reda Bahi Slaoui, and Sébastien Toth. Estimating risk and uncertainty in deep reinforcement learning. *arXiv preprint arXiv:1905.09638*, 2019. [3](#)
- [13] Emily Denton and Rob Fergus. Stochastic video generation with a learned prior. In *International Conference on Machine Learning*, pages 1174–1183, 2018. [1](#), [2](#), [3](#), [4](#), [6](#), [7](#), [8](#), [12](#)
- [14] Emily L Denton et al. Unsupervised learning of disentangled representations from video. In *Advances in neural information processing systems*, pages 4414–4423, 2017. [3](#)
- [15] Frederik Ebert, Chelsea Finn, Alex X Lee, and Sergey Levine. Self-supervised visual planning with temporal skip connections. *arXiv preprint arXiv:1710.05268*, 2017. [1](#), [2](#), [6](#)
- [16] Mikhail Figurnov, Shakir Mohamed, and Andriy Mnih. Implicit reparameterization gradients. In *Advances in Neural Information Processing Systems*, pages 441–452, 2018. [5](#)
- [17] Chelsea Finn, Ian Goodfellow, and Sergey Levine. Unsupervised learning for physical interaction through video prediction. In *Advances in neural information processing systems*, pages 64–72, 2016. [2](#)
- [18] Marco Fraccaro, Søren Kaae Sønderby, Ulrich Paquet, and Ole Winther. Sequential neural models with stochastic layers. In *Advances in neural information processing systems*, pages 2199–2207, 2016. [2](#)
- [19] Katerina Fragkiadaki, Sergey Levine, Panna Felsen, and Jitendra Malik. Recurrent network models for human dynamics. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4346–4354, 2015. [1](#), [2](#)
- [20] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014. [1](#), [2](#), [6](#)
- [21] Vincent Le Guen and Nicolas Thome. Disentangling physical dynamics from unknown factors for unsupervised video prediction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11474–11484, 2020. [3](#)
- [22] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville. Improved training of wasserstein gans. In *Advances in neural information processing systems*, pages 5767–5777, 2017. [6](#)
- [23] Shir Gur, Sagie Benaim, and Lior Wolf. Hierarchical patch vae-gan: Generating diverse videos from a single sample. *Advances in Neural Information Processing Systems*, 33, 2020. [6](#)
- [24] Ali Harakeh, Michael Smart, and Steven L Waslander. Bayesod: A bayesian approach for uncertainty estimation in deep object detectors. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 87–93. IEEE, 2020. [3](#)
- [25] Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual concepts with a constrained variational framework. *ICLR*, 2(5):6, 2017. [5](#)
- [26] Yung-Han Ho, Chuan-Yuan Cho, Wen-Hsiao Peng, and Guo-Lun Jin. Sme-net: Sparse motion estimation for parametric video prediction through reinforcement learning. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 10462–10470, 2019. [3](#)
- [27] Jun-Ting Hsieh, Bingbin Liu, De-An Huang, Li F Fei-Fei, and Juan Carlos Niebles. Learning to decompose and disentangle representations for video prediction. In *Advances in Neural Information Processing Systems*, pages 517–526, 2018. [3](#), [6](#), [7](#)
- [28] Ashesh Jain, Hema S Koppula, Bharad Raghavan, Shane Soh, and Ashutosh Saxena. Car that knows before you do:

- Anticipating maneuvers via learning temporal driving models. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3182–3190, 2015. 1, 2
- [29] Beibei Jin, Yu Hu, Qiankun Tang, Jingyu Niu, Zhiping Shi, Yinhe Han, and Xiaowei Li. Exploring spatial-temporal multi-frequency analysis for high-fidelity and temporal-consistency video prediction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4554–4563, 2020. 1, 2, 3
- [30] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4401–4410, 2019. 1
- [31] Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer vision? In *Advances in neural information processing systems*, pages 5574–5584, 2017. 2
- [32] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. 2015. 6
- [33] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. 2014. 3, 4, 5, 12
- [34] David A Knowles. Stochastic gradient variational bayes for gamma approximating distributions. *arXiv preprint arXiv:1509.01631*, 2015. 5
- [35] Rahul G Krishnan, Uri Shalit, and David Sontag. Structured inference networks for nonlinear state space models. In *Thirty-first aaai conference on artificial intelligence*, 2017. 2
- [36] Abhinav Kumar, Tim K Marks, Wenxuan Mou, Ye Wang, Michael Jones, Anoop Cherian, Toshiaki Koike-Akino, Xiaoming Liu, and Chen Feng. Luvli face alignment: Estimating landmarks’ location, uncertainty, and visibility likelihood. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8236–8246, 2020. 3
- [37] Yong-Hoon Kwon and Min-Gyu Park. Predicting future frames using retrospective cycle gan. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1811–1820, 2019. 1, 2, 6
- [38] Minhyeok Lee and Junhee Seok. Estimation with uncertainty via conditional generative adversarial networks. *arXiv preprint arXiv:2007.00334*, 2020. 3
- [39] Xiaodan Liang, Lisa Lee, Wei Dai, and Eric P Xing. Dual motion gan for future-flow embedded video prediction. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1744–1752, 2017. 1, 2, 6
- [40] Wen Liu, Weixin Luo, Dongze Lian, and Shenghua Gao. Future frame prediction for anomaly detection—a new baseline. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6536–6545, 2018. 3
- [41] William Lotter, Gabriel Kreiman, and David Cox. Deep predictive coding networks for video prediction and unsupervised learning. *International Conference on Learning Representations*, 2017. 2
- [42] Andrey Malinin and Mark Gales. Predictive uncertainty estimation via prior networks. In *Advances in Neural Information Processing Systems*, pages 7047–7058, 2018. 3
- [43] Rowan McAllister, Gregory Kahn, Jeff Clune, and Sergey Levine. Robustness to out-of-distribution inputs via task-aware generative uncertainty. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 2083–2089. IEEE, 2019. 3
- [44] MarcAurelio Ranzato, Arthur Szlam, Joan Bruna, Michael Mathieu, Ronan Collobert, and Sumit Chopra. Video (language) modeling: a baseline for generative models of natural videos. *arXiv preprint arXiv:1412.6604*, 2014. 2
- [45] Christian P Robert. Simulation of truncated normal variables. *Statistics and computing*, 5(2):121–125, 1995. 6
- [46] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. 12
- [47] Francisco R Ruiz, Michalis Titsias RC AUEB, and David Blei. The generalized reparameterization gradient. In *Advances in neural information processing systems*, pages 460–468, 2016. 5
- [48] Christian Schuldt, Ivan Laptev, and Barbara Caputo. Recognizing human actions: a local svm approach. In *Proceedings of the 17th International Conference on Pattern Recognition, 2004. ICPR 2004.*, volume 3, pages 32–36. IEEE, 2004. 2, 6
- [49] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. UCF101: A dataset of 101 human actions classes from videos in the wild. *CoRR*, abs/1212.0402, 2012. 2, 6, 13
- [50] Nitish Srivastava, Elman Mansimov, and Ruslan Salakhudinov. Unsupervised learning of video representations using lstms. In *International conference on machine learning*, pages 843–852, 2015. 2
- [51] Chen Sun, Abhinav Shrivastava, Carl Vondrick, Rahul Sukthankar, Kevin Murphy, and Cordelia Schmid. Relational action forecasting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 273–283, 2019. 1
- [52] Natasa Tagasovska and David Lopez-Paz. Single-model uncertainties for deep learning. In *Advances in Neural Information Processing Systems*, pages 6417–6428, 2019. 3
- [53] Jakub Tomczak and Max Welling. Vae with a vampprior. In *International Conference on Artificial Intelligence and Statistics*, pages 1214–1223, 2018. 3
- [54] Arash Vahdat and Jan Kautz. Nvae: A deep hierarchical variational autoencoder. *Advances in Neural Information Processing Systems*, 33, 2020. 3
- [55] Ruben Villegas, Jimei Yang, Seunghoon Hong, Xunyu Lin, and Honglak Lee. Decomposing motion and content for natural video sequence prediction. *Proceedings of the International Conference on Learning Representations*, 2017. 3
- [56] Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. Generating videos with scene dynamics. In *Advances in neural information processing systems*, pages 613–621, 2016. 1
- [57] Chandrasekar Vuppapapati, Anitha Ilapakurti, Sharat Kedari, Rajasekar Vuppapapati, Jayashankar Vuppapapati, and Santosh Kedari. Human ai symbiosis: The role of artificial intelligence in stratifying high-risk outpatient senior citizen fall events in a non-connected environments. In *International Conference on Intelligent Human Systems Integration*, pages 325–332. Springer, 2020. 1, 2

- [58] Jacob Walker, Kenneth Marino, Abhinav Gupta, and Martial Hebert. The pose knows: Video forecasting by generating pose futures. In *Proceedings of the IEEE international conference on computer vision*, pages 3332–3341, 2017. 1
- [59] Gang Wang, Bo Li, Yongfei Zhang, and Jinhui Yang. Background modeling and referencing for moving cameras-captured surveillance video coding in hev. *IEEE Transactions on Multimedia*, 20(11):2921–2934, 2018. 1
- [60] Max Welling and Yee W Teh. Bayesian learning via stochastic gradient langevin dynamics. In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pages 681–688, 2011. 3
- [61] Yue Wu, Rongrong Gao, Jaesik Park, and Qifeng Chen. Future video synthesis with object motion prediction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5539–5548, 2020. 3
- [62] SHI Xingjian, Zhoung Chen, Hao Wang, Dit-Yan Yeung, Wai-Kin Wong, and Wang-chun Woo. Convolutional lstm network: A machine learning approach for precipitation nowcasting. In *Advances in neural information processing systems*, pages 802–810, 2015. 2
- [63] Yufei Ye, Maneesh Singh, Abhinav Gupta, and Shubham Tulsiani. Compositional video prediction. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 10353–10362, 2019. 3
- [64] Wei Yu, Yichao Lu, Steve Easterbrook, and Sanja Fidler. Efficient and information-preserving future frame prediction and beyond. In *International Conference on Learning Representations*, 2019. 2
- [65] 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. 6

## A. Uncertainty Visualizations

Figure 7 visualizes the scaled uncertainty values against the visual frames, across the SMMNIST and BAIR Push datasets, each trained with 2000 samples. See the caption of the figure for more details.

## B. Alternative Formulations

Is NUQ the best formulation that one could have for quantifying uncertainty within a stochastic prediction model? In this section, we propose several alternatives and empirically evaluate them against the results we obtained using our formulation of NUQ, as an answer to this interesting research question.

### B.1. Alternatives

**Alternative Priors:** In this variant, we replace our empirical gamma hyperprior,  $p(s_t)$ , with a half-normal distribution with location set to 0 and scale set to 1, and in another variant we use the Uniform(0, 1) distribution as the hyperprior.

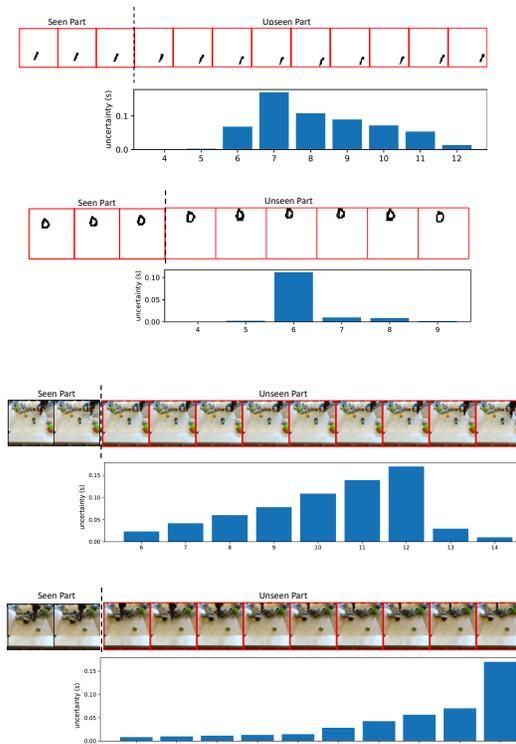


Figure 7. Top two rows: Evolution of scaled uncertainty on the SMMNIST dataset (with NUQ trained with 2,000 training samples on the SMMNIST Dataset) against time-steps. The plot shows the increase in uncertainty co-occurs with the bounce of the digit against the boundary. Bottom two rows: Scaled Uncertainty against time-steps on BAIR Push (with NUQ trained with 2000 training samples on the BAIR Push Dataset), showing that uncertainty co-occurs with the occlusion of the robot arm.

**Using Mahalanobis Distance:** In this variant, we use the frame decoder,  $p_\theta(\cdot)$ , to produce a diagonal covariance instead of producing the parameters of the gamma prior. Specifically, the output layer of the decoder now predicts both the future frame and the diagonal elements of  $\Sigma_t^z$ , where we assume  $\Sigma_t^z$  is  $n \times n$ , decoded frames are size  $d \times d$ , and  $D = d^2$ . Generating an estimate of the output covariance matrix thus implies predicting the  $D$  terms along the diagonal of this matrix. We then reshape these  $D$  terms to  $d \times d$  to match with the pixel resolution of the frames. We then use this uncertainty (precision) to weigh the MSE at a pixel level (instead of the precision  $b_t$ ). This uncertainty is visualized for a sequence in Figure 8.



Figure 8. Visualization of pixel-wise uncertainty obtained by estimating the variance of the output, directly from the decoder for the SMMNIST Dataset, trained with 2,000 training samples.

**Directly Estimate precision from latent prior:** In this formulation, we repurpose the variance encoder  $\zeta_\lambda$ , to emit the variance to the MSE,  $b_t$ , directly. This is in contrast with the architecture of the variance encoder in NUQ, where it is used to estimate the sufficient statistics of the truncated normal distribution,  $\alpha_t^{\tilde{s}}$  and  $\beta_t^{\tilde{s}}$ . We essentially replace the final hidden layer of the network with a single neuron with sigmoid activation in order to realize this setting.

## B.2. Alternatives – Results

In Table 4, we provide comparisons of the above alternative formulations on the SMMNIST dataset, trained with 2,000 training samples. From the first row, we see that the Uniform[0, 1] distribution variant under-performs compared to using the gamma distribution as a hyper-prior, as in NUQ. We surmise that this is due to these distributions being more spread out over the probability space, as a result of which they often sample  $s_t$ 's which do not match the true underlying distribution. This results in the MSE term in the loss function, being overly weighed when it should not be and vice-versa.

Results for our other alternative, to compute the Mahalanobis-type precision matrix directly from the frame decoder, is provided in the second row in Table 4; its performance is similar to the other alternatives. We also attempted to directly estimate  $s_t$  from the variance of the latent space prior  $\Sigma_t^z$ . The results for this setting are shown in the third row in Table 4. However, this setting performs poorly suggesting that a deterministic mapping of the  $\Sigma_t^z$  to  $s_t$  is not ideal, perhaps because the difference in the uncertainty distribution in the latent space and in the output is not accurately modeled this way. Overall, the results in the table clearly show that our proposed formulation of the model yields the best empirical performance, nonetheless some other formulations to our model seem promising.

## C. Architectural Details

In this section, we elaborate on some of the architectural design choices that we made while implementing NUQ. Our primary objective while designing the architectural framework of NUQ was to ensure that our network's generation capacity remained similar to the state-of-the-art baselines, such as Denton and Fergus [13], such that all gains obtained by our framework, could be attributed to modeling the prediction uncertainty.

### C.1. Frame Encoder

Our frame encoder consists of a hierarchical stack of 2d-convolution filters. For  $48 \times 48$  inputs, we design a 4-layer network. The first layer consists of 64,  $4 \times 4$  2d-convolutional filters with stride 2 and padding 1, which are followed by 2d-BatchNorm and LeakyReLU non-linearity. In every subsequent layers, we keep doubling the number of

filters. For  $64 \times 64$  inputs, we adapt this network to make it a 5-layer one.

### C.2. LSTMs

All LSTM modules in our framework, including the sequence discriminator, have a single hidden layer with 256-d hidden states, except for the LSTM in the frame decoder  $p_\theta(\cdot)$ , which has 2 hidden layers, each of 256-d.

### C.3. Frame Decoder

We design the frame decoder in congruence with the frame encoder, so as to permit skip connections between them, in a U-Net style network [46]. Therefore, our frame decoder obeys a similar architecture akin to the frame encoder, except the 2d-convolution filters are now replaced with 2d-deconvolution filters and the number of filters in each layer is doubled (in order to accommodate the skip connection).

### C.4. Variance Encoder

Our variance encoder,  $\zeta_\lambda(\cdot)$ , is a 2-layer multi-layer perceptron, with LeakyReLU activations, which ultimately produces the sufficient statistics of the truncated normal distribution governing the posterior in the latent space.

## D. Derivations

In this section, we present the derivations of Eq. 9 and Eq. 11 in the paper. We derive Eq. 9, along the lines of the variational lower bound derivation in Kingma and Welling [33]:

$$\begin{aligned} \log p(\mathbf{x}_t | b_t, \mathbf{x}_{1:t-1}) &= \log p(\mathbf{x}_t | b_t, \mathbf{x}_{1:t-1}) \cdot \int_{\mathbf{z}_t} q_\phi(\mathbf{z}_t | \mathbf{x}_{1:t}) d\mathbf{z}_t \\ &= \int_{\mathbf{z}_t} q_\phi(\mathbf{z}_t | \mathbf{x}_{1:t}) \log \frac{p(\mathbf{x}_t, \mathbf{z}_t | b_t, \mathbf{x}_{1:t-1})}{p(\mathbf{z}_t | \mathbf{x}_{1:t})} d\mathbf{z}_t \\ &= \int_{\mathbf{z}_t} q_\phi(\mathbf{z}_t | \mathbf{x}_{1:t}) \log \frac{p(\mathbf{x}_t, \mathbf{z}_t | b_t, \mathbf{x}_{1:t-1})}{q_\phi(\mathbf{z}_t | \mathbf{x}_{1:t})} d\mathbf{z}_t \\ &\quad + \int_{\mathbf{z}_t} q_\phi(\mathbf{z}_t | \mathbf{x}_{1:t}) \log \frac{q_\phi(\mathbf{z}_t | \mathbf{x}_{1:t})}{p(\mathbf{z}_t | \mathbf{x}_{1:t})} d\mathbf{z}_t \end{aligned} \tag{15}$$

The second term in the above equation is essentially a KL-Divergence, which is non-negative. We therefore have:

$$\begin{aligned} \log p(\mathbf{x}_t | b_t, \mathbf{x}_{1:t-1}) &\geq \int_{\mathbf{z}_t} q_\phi(\mathbf{z}_t | \mathbf{x}_{1:t}) \log \frac{p(\mathbf{x}_t, \mathbf{z}_t | b_t, \mathbf{x}_{1:t-1})}{q_\phi(\mathbf{z}_t | \mathbf{x}_{1:t})} d\mathbf{z}_t \\ &= \int_{\mathbf{z}_t} q_\phi(\mathbf{z}_t | \mathbf{x}_{1:t}) \log \frac{p(\mathbf{x}_t | \mathbf{x}_{1:t-1}, \mathbf{z}_t, b_t) p(\mathbf{z}_t | \mathbf{x}_{1:t-1})}{q_\phi(\mathbf{z}_t | \mathbf{x}_{1:t})} d\mathbf{z}_t \end{aligned} \tag{16}$$

This yields Eq. 9, when the expression inside the log is split into two, with the first term amounting to the expectation term in Eq. 9, while the second one resulting in the KL-term.

Table 4. SSIM, PSNR, and LPIPS scores on the SMMNIST test set after @1, @5, and @Convergence (C) (upto 150) epochs of training with alternative formulations of our model using 2,000 training samples. [Key: Best results in **bold**].

Dataset: SMMNIST	SSIM ↑			PSNR ↑			LPIPS ↓		
	@1	@5	@C	@1	@5	@C	@1	@5	@C
$p(s_t) \sim \text{Uniform}[0, 1]$	0.8173	0.8374	0.8523	17.6	17.95	18.06	0.3442	0.3038	0.198
Estimate $b_t$ from the decoder $p_\theta(\cdot)$	0.7627	0.7628	0.7828	17.54	17.55	17.55	0.3463	0.3259	0.2225
Estimate $b_t$ w/o variance enc-dec	0.7450	0.7454	0.7648	16.22	16.53	16.78	0.3469	0.3263	0.2328
NUQ (Ours)	<b>0.8686</b>	<b>0.8638</b>	<b>0.8948</b>	<b>17.76</b>	<b>18.13</b>	<b>18.14</b>	<b>0.3087</b>	<b>0.2836</b>	<b>0.1803</b>

Our NUQ framework is essentially, a hierarchical variational encoder-decoder network, where the second level of the hierarchy is described by Eq. 11. Derivation for Eq. 11, thus proceeds analogously to Eq. 9, as follows:

$$\begin{aligned}
 \log p(b_t | \mathbf{x}_{1:t-1}) &= \log p(b_t | \mathbf{x}_{1:t-1}) \cdot \int_{s_t} q_\lambda(s_t | \mathbf{x}_{1:t-1}) ds_t \\
 &= \int_{s_t} q_\lambda(s_t | \mathbf{x}_{1:t-1}) \log \frac{p(b_t, s_t | \mathbf{x}_{1:t-1})}{p(s_t | b_t, \mathbf{x}_{1:t-1})} ds_t \\
 &= \int_{s_t} q_\phi(s_t | \mathbf{x}_{1:t-1}) \log \frac{p(b_t, s_t | \mathbf{x}_{1:t-1})}{q_\lambda(s_t | \mathbf{x}_{1:t-1})} ds_t \\
 &+ \int_{s_t} q_\lambda(s_t | \mathbf{x}_{1:t-1}) \log \frac{q_\lambda(s_t | \mathbf{x}_{1:t-1})}{p(s_t | b_t, \mathbf{x}_{1:t-1})} ds_t
 \end{aligned}
 \tag{17}$$

Like before, the second term in the above equation is essentially a KL-Divergence, which is non-negative. We therefore have:

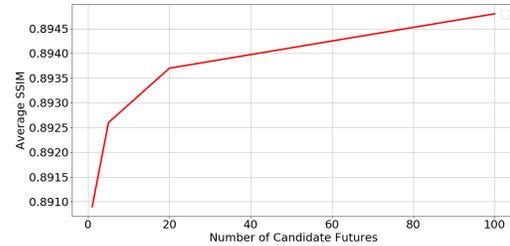
$$\begin{aligned}
 \log p(b_t | \mathbf{x}_{1:t-1}) &\geq \int_{s_t} q_\lambda(s_t | \mathbf{x}_{1:t-1}) \log \frac{p(b_t, s_t | \mathbf{x}_{1:t-1})}{q_\lambda(s_t | \mathbf{x}_{1:t-1})} ds_t \\
 &= \int_{s_t} q_\lambda(s_t | \mathbf{x}_{1:t-1}) \log \frac{p(b_t | \mathbf{x}_{1:t-1}, s_t) p(s_t)}{q_\lambda(s_t | \mathbf{x}_{1:t-1})} ds_t
 \end{aligned}
 \tag{18}$$

When the expression inside the log is split into two, the first term results in the expectation term in Eq. 11, while the second one amounts to the KL-term.

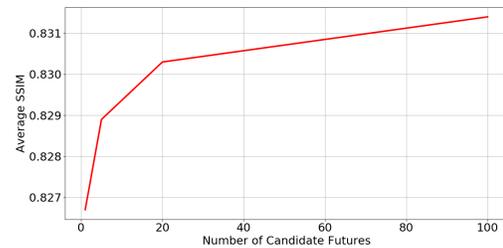
In this section, we present model performances of NUQ versus competing baselines on the UCF-101 dataset [49] – a dataset of videos, resized to  $64 \times 64$  containing 101 common human action categories (such as pushups, cricket shot, etc.), spanning both indoor and outdoor activities. The test set for this dataset consists of 1,895 videos. In order to conduct experiments in this setting, we train all models by showing them 5 context frames and task them to predict the next 15. The results showcase the dominance of NUQ over competing methods on this challenging dataset as well.

## E. Quantitative Evaluation of Diversity

In order to analyze the extent of diversity in the generated frames of our model, we first resort to quantitative evaluation. In Figures 9(a), 9(b), we plot the average SSIM scores (over time steps) against the number of generated future



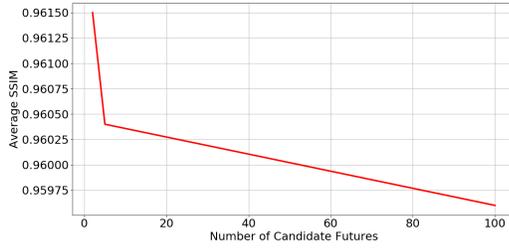
(a) SMMNIST - SSIM



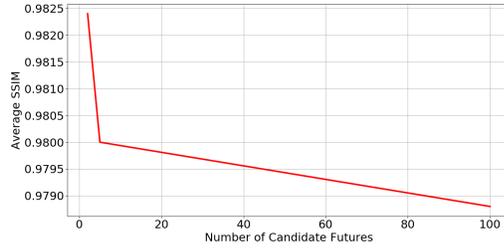
(b) BAIR Push - SSIM

Figure 9. Diversity in Generated Futures: Evaluation of diversity in the generation using SSIM on: (a) SMMNIST, (b) BAIR-Push for increasing number of candidate futures, computed by comparing against the ground truth (higher the better). We used 2000 samples for training NUQ for both datasets.

candidates per time-step for each of the three datasets. For purposes of these plots, the SSIM is computed between the generated samples and the ground-truth. The monotonically increasing curve, in these figures, suggests straightforwardly, that sampling more future per time step helps in better generation, resulting from the synthesis of more accurate samples - indicating the model’s diversity. In Figures 10(a), 10(b), we plot the average SSIM and PSNR scores between every pair of candidates generated in each time step, against the number of futures. These plots decrease monotonically, suggesting greater difference (i.e. diversity) between the generated frames as the number of sampled futures goes up. See the figure caption for more details.



(a) SMMNIST-Intra



(b) BAIR-Intra

Figure 10. Diversity in Generated Futures: Evaluation of diversity in the generation: (a,b) shows diversity in the generated futures by comparing intra-SSIM distances between all the futures, at a given time step, and computing the average (lower the better), for SMMNIST and BAIR Push respectively. For each of these datasets NUQ was trained with 2000 samples.

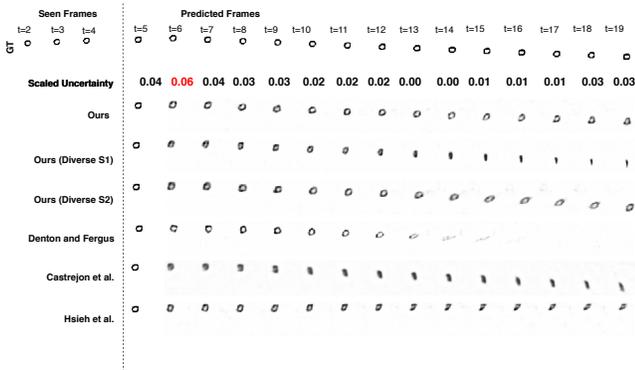


Figure 11. Visualization of generations by our method versus competing baselines on the SMMNIST Dataset, trained with 2,000 training samples. Further, diverse generations by our method are also shown. Note scaled uncertainty higher than 0.05 is shown in red.

## F. Qualitative Results

We next present visualizations of frames generated using NUQ vis-à-vis competing baselines, on the SMMNIST, BAIR push, UCF-101, and KTH Action datasets. Also shown are diverse frame generations by NUQ for each of these datasets.

The results in Figures 11, 14, 15, 16, 17, 18, 19, 20, 21 show qualitative generation results on the SMMNIST dataset,

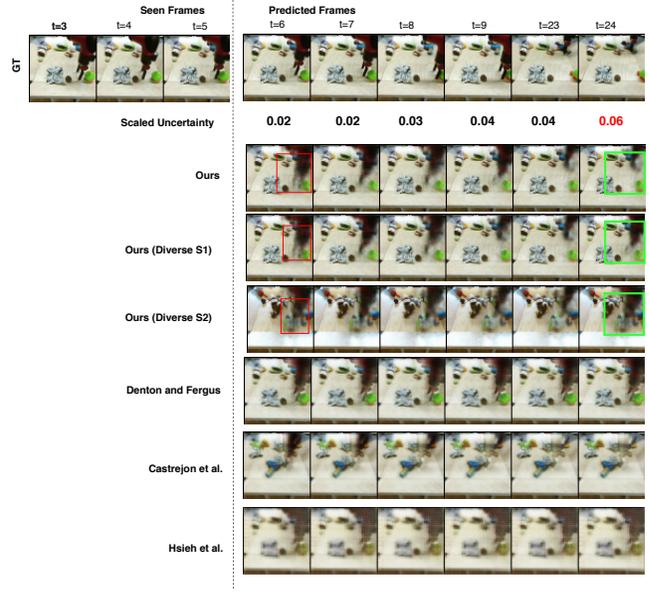


Figure 12. Visualization of generations by our method versus competing baselines on the BAIR Robot Push Dataset, trained with 2,000 training samples. Further, diverse generations by our method are also shown. High motion regions are indicated by a red bounding box, while spatial regions exhibiting high diversity are shown by a green bounding box. Note scaled uncertainty higher than 0.05 is shown in red.

trained with 2000 samples. Besides the superior quality of the results generated by our method, we note that for some cases such as in Figures 17, 18, 19 the prediction of the baseline method simply disappears. We surmise that this is due to their inability to learn the motion dynamics of the digit well, in uncertain environments. In particular, if the motion is not aptly learnt, then the model often gets penalized heavily for inaccurately placing a digit (via the MSE loss), since this results in a high pixel-wise error. In such a scenario, a model might prefer to not display the digit at all. However, high stochasticity in the data may not suit this well and as a result might hurt the generalization. By intelligently down-weighting the MSE, we circumvent this problem.

We also see in Figures 12, 22, 26, 23, 27, 28, 29, 30, 24, 25 the performance of different competing methods versus NUQ on the BAIR push dataset, trained with 2000 samples. The figures reveal that our method captures the motion of the robot arm, reasonably well, compared to competing methods.

Figures 13, 31, 32 present sample generation results by our method versus competing baselines on the KTH Action dataset. From the figures, we see that while all of the methods do a reasonable job of modeling the appearance of the person, nonetheless the competing methods fail to capture the motion dynamics well.

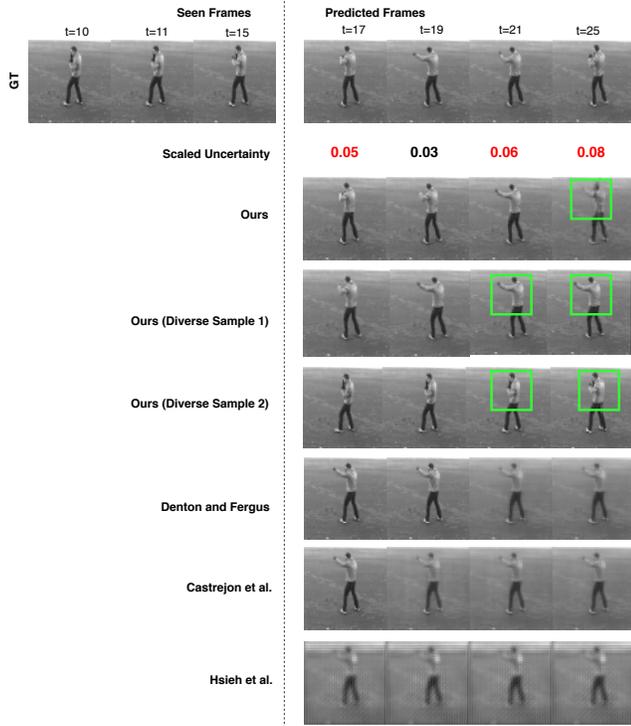


Figure 13. Visualization of generations by our method versus competing baselines on the KTH Action Dataset, trained with the full training data of 1,911 training samples. Further, diverse generations by our method are also shown. Spatial regions exhibiting high diversity are shown by a green bounding box. Note scaled uncertainty higher than 0.05 is shown in red.

Such trends extend into the UCF-101 dataset as well. The results for this dataset are shown in Figures 33, 34.

Moreover, in some of the aforementioned figures (such as Figures 14, 15, 16, 22, 12, 25, 31, 32, 33, 34 diverse sample generations by NUQ is also shown.

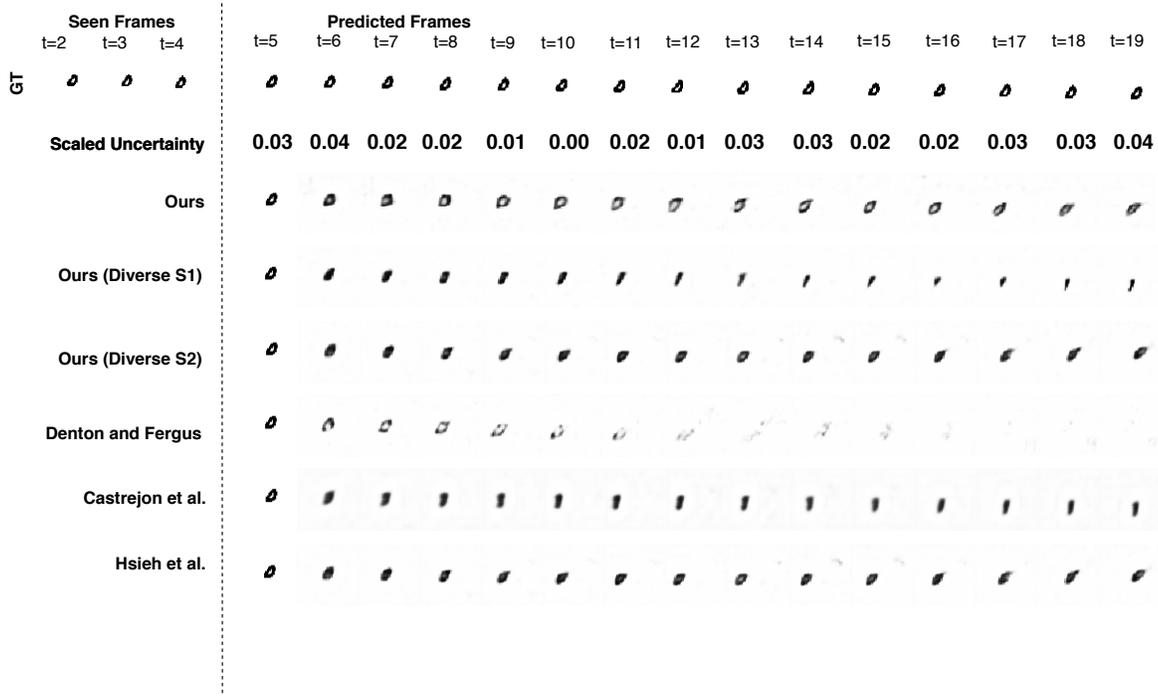


Figure 14. Visualization of generations by our method versus competing baselines on the SMMNIST Dataset, trained with 2,000 training samples. Further, diverse generations by our method are also shown. Note scaled uncertainty higher than 0.05 is shown in red.

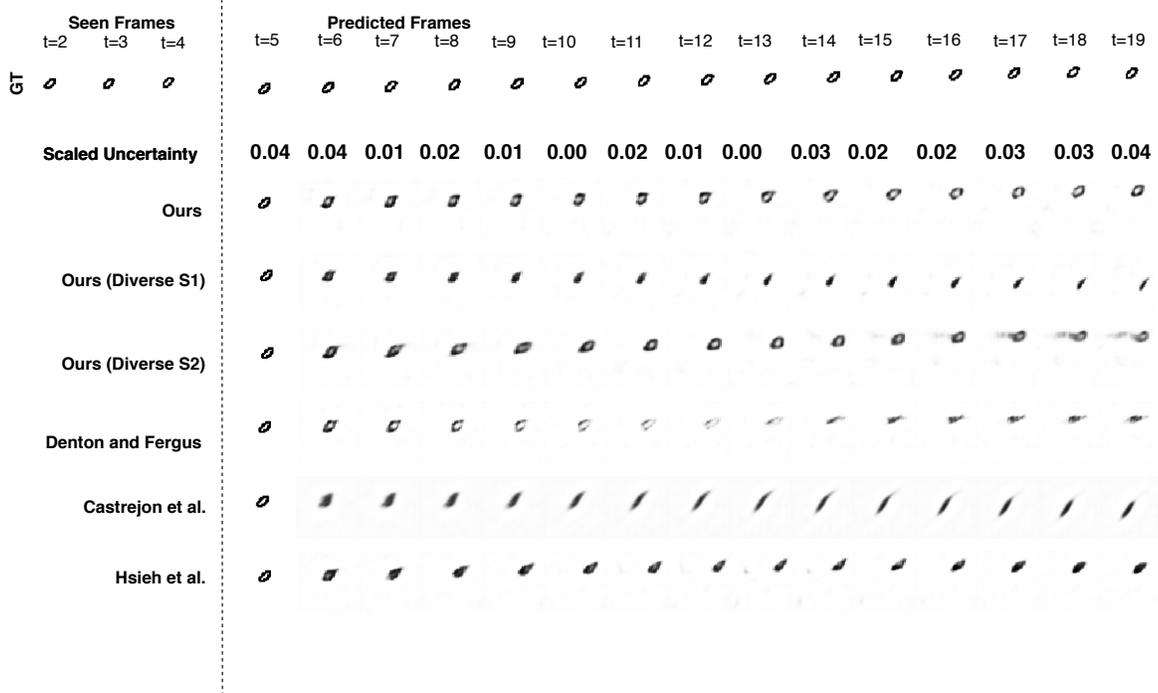


Figure 15. Visualization of generations by our method versus competing baselines on the SMMNIST Dataset, trained with 2,000 training samples. Further, diverse generations by our method are also shown. Note scaled uncertainty higher than 0.05 is shown in red.

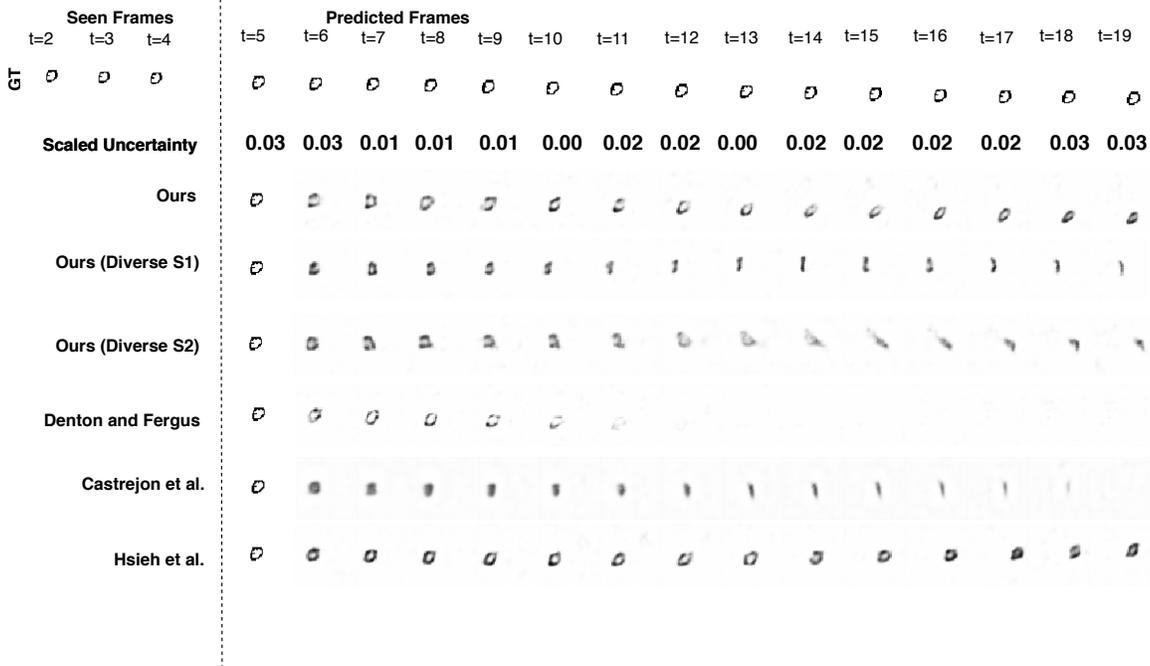


Figure 16. Visualization of generations by our method versus competing baselines on the SMMNIST Dataset, trained with 2,000 training samples. Further, diverse generations by our method are also shown. Note scaled uncertainty higher than 0.05 is shown in red.

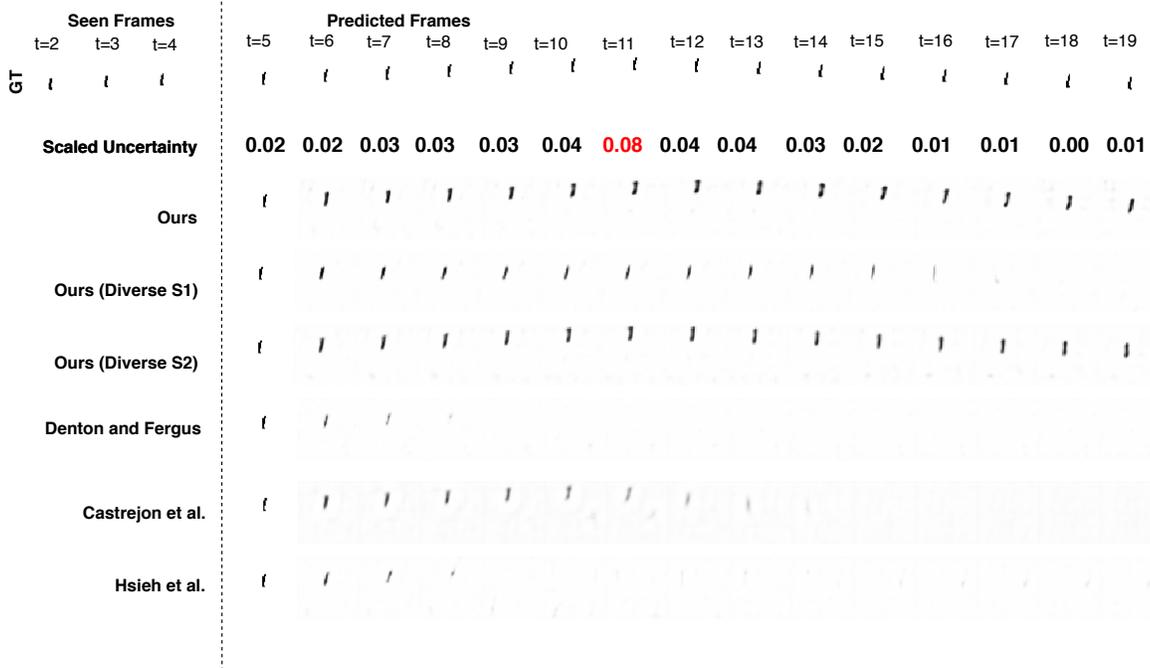


Figure 17. Visualization of generations by our method versus competing baselines on the SMMNIST Dataset, trained with 2,000 training samples. Further, diverse generations by our method are also shown. Note scaled uncertainty higher than 0.05 is shown in red.

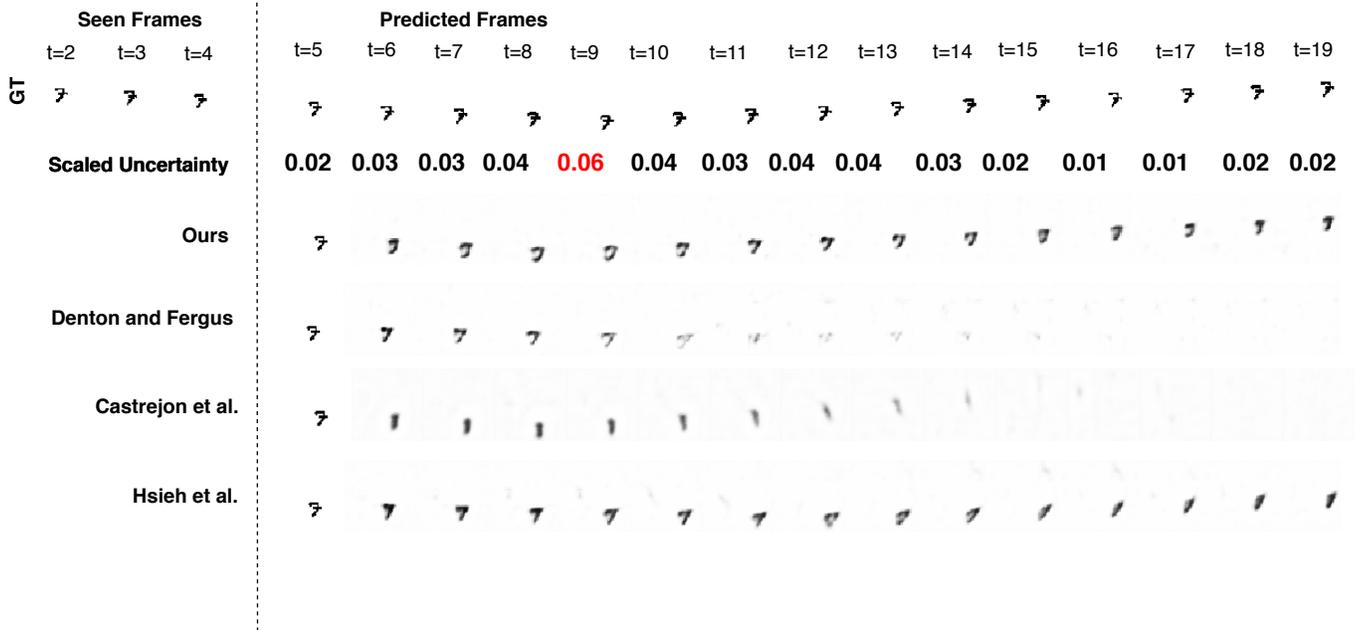


Figure 18. Visualization of generations by our method versus competing baselines on the SMMNIST Dataset, trained with 2,000 training samples. Note scaled uncertainty higher than 0.05 is shown in red.

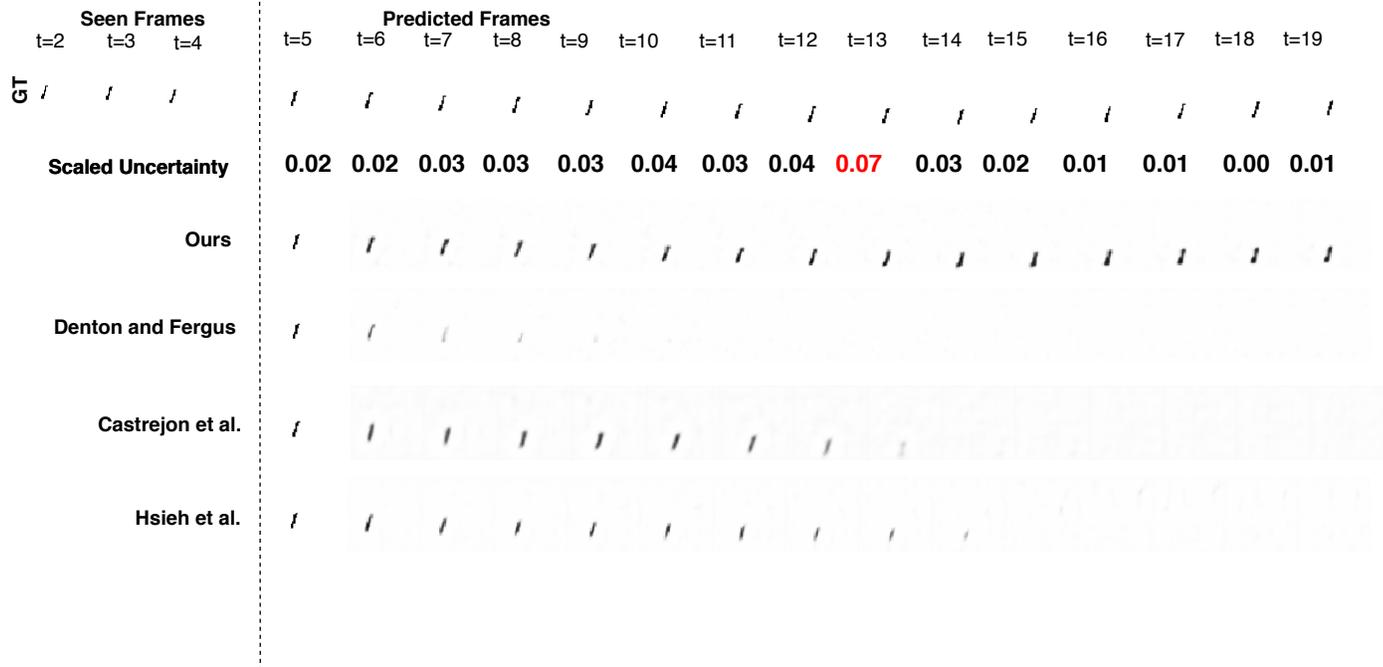


Figure 19. Visualization of generations by our method versus competing baselines on the SMMNIST Dataset, trained with 2,000 training samples. Note scaled uncertainty higher than 0.05 is shown in red.

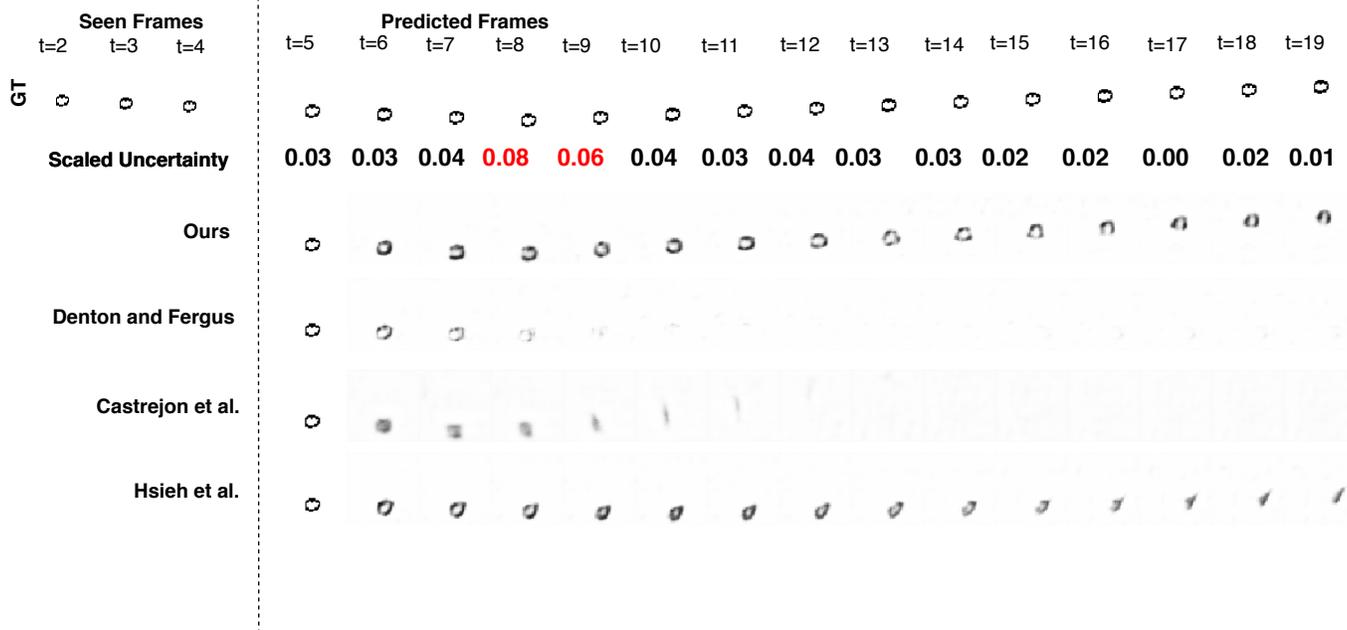


Figure 20. Visualization of generations by our method versus competing baselines on the SMMNIST Dataset, trained with 2,000 training samples. Note scaled uncertainty higher than 0.05 is shown in red.

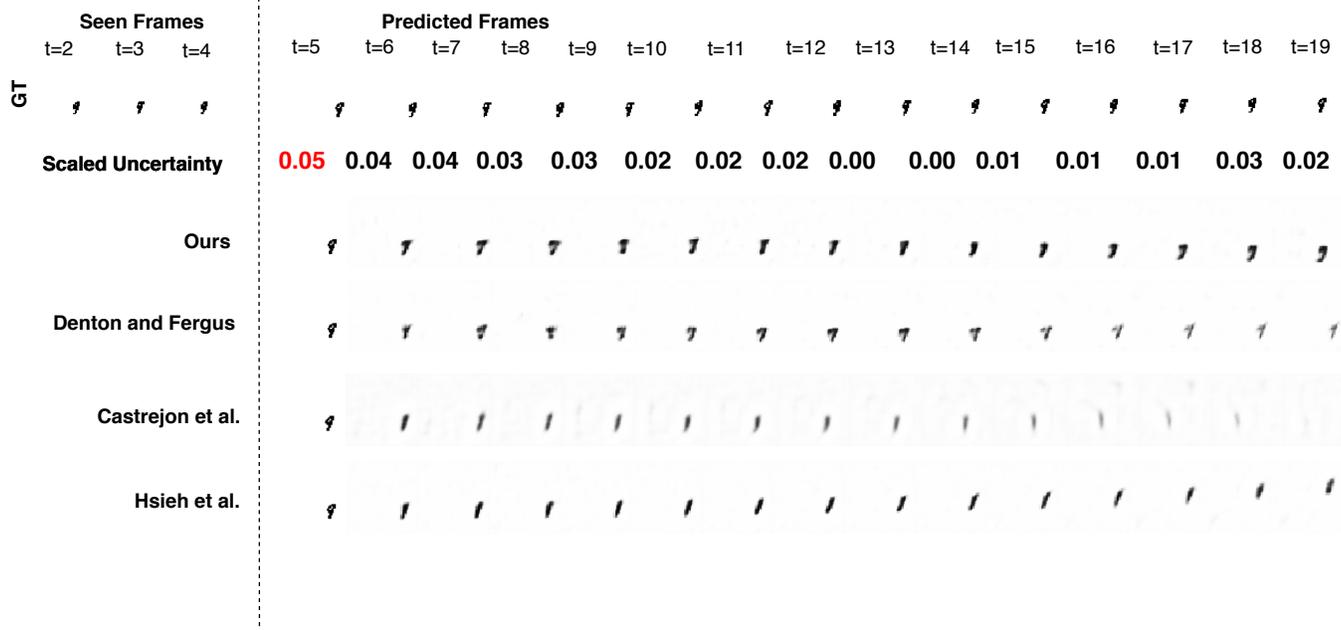


Figure 21. Visualization of generations by our method versus competing baselines on the SMMNIST Dataset, trained with 2,000 training samples. Note scaled uncertainty higher than 0.05 is shown in red.

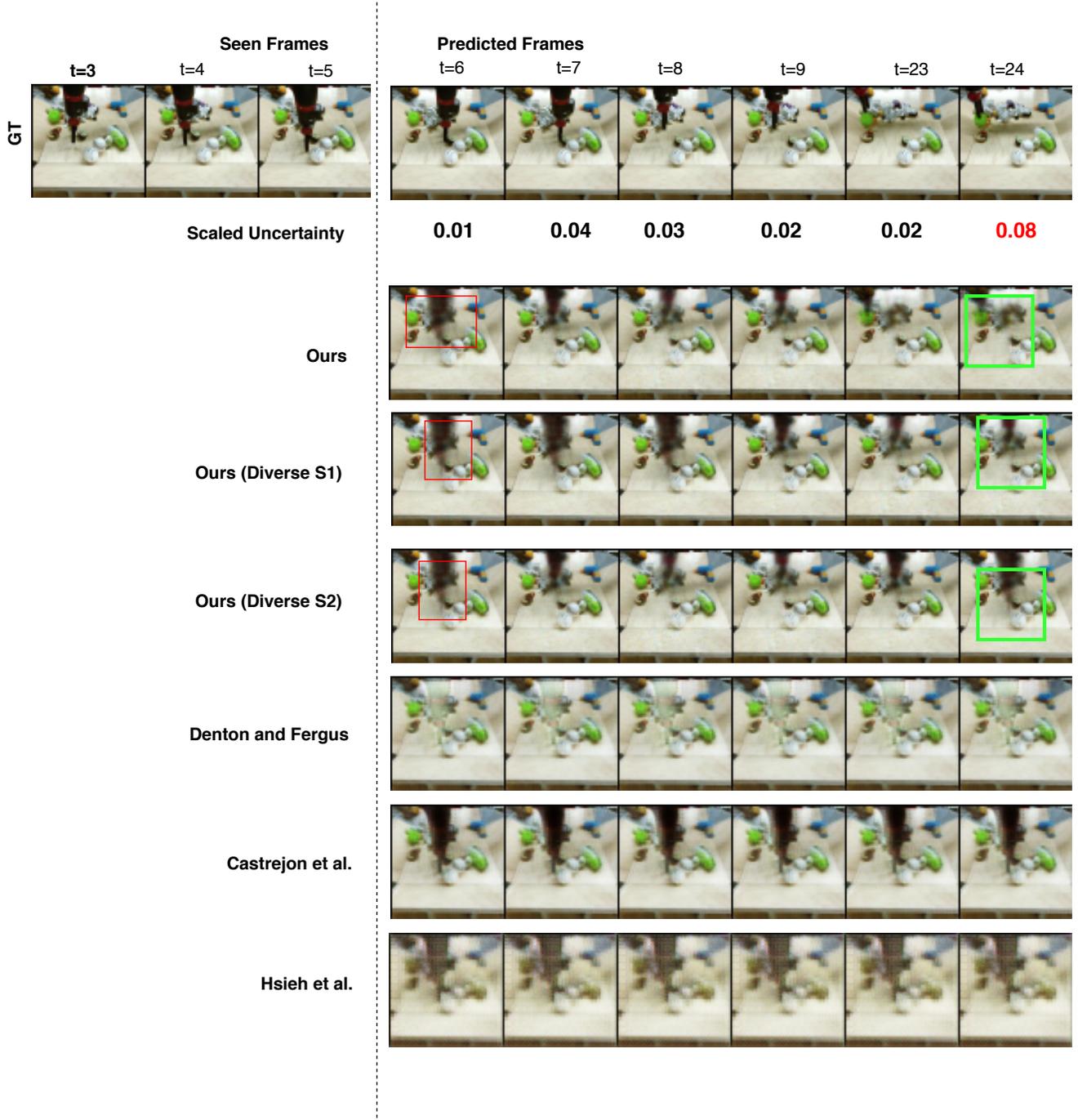


Figure 22. Visualization of generations by our method versus competing baselines on the BAIR Robot Push Dataset, trained with 2,000 training samples. Further, diverse generations by our method are also shown. High motion regions are indicated by a red bounding box, while spatial regions exhibiting high diversity are shown by a green bounding box. Note scaled uncertainty higher than 0.05 is shown in red.

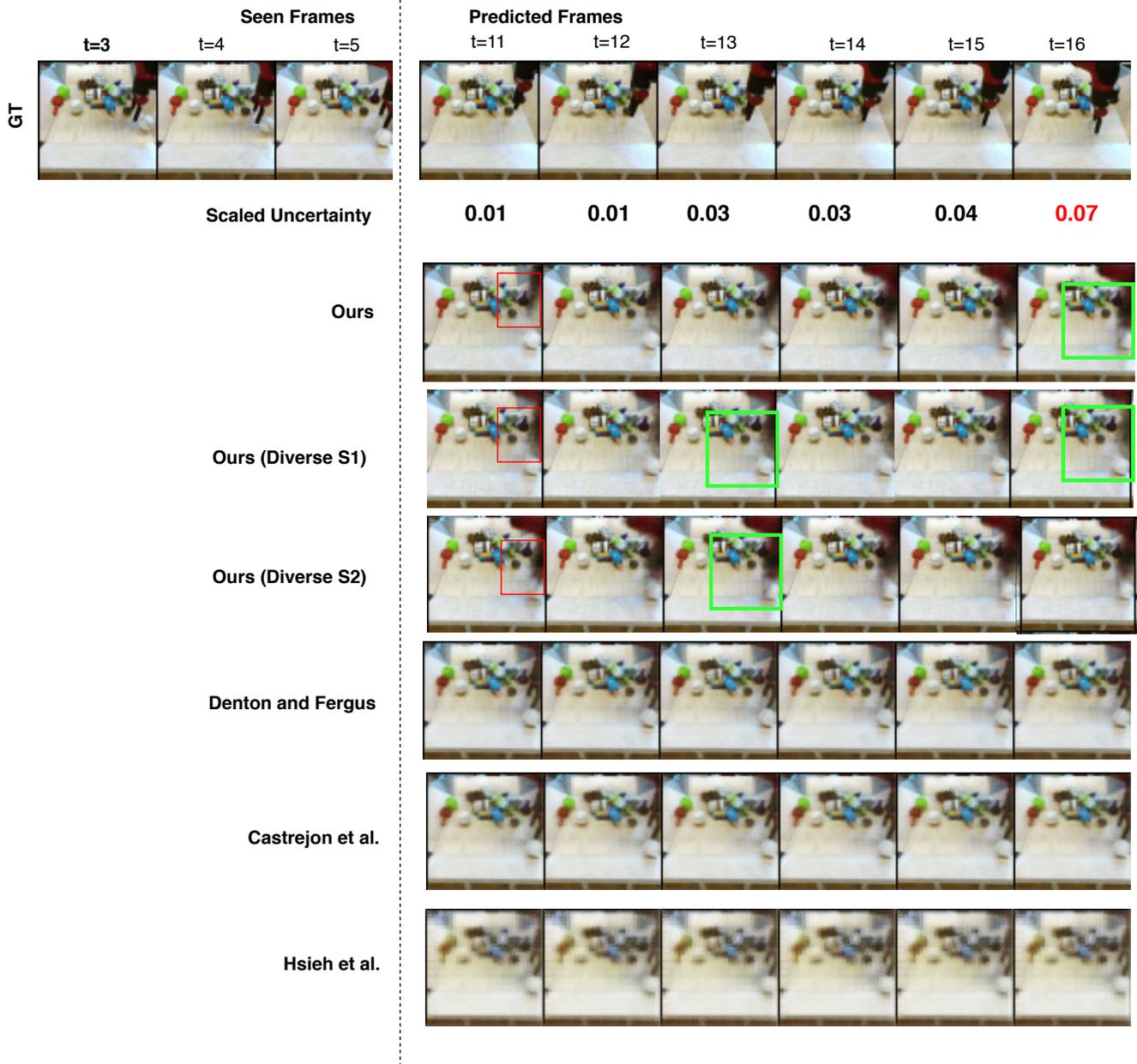


Figure 23. Visualization of generations by our method versus competing baselines on the BAIR Robot Push Dataset, trained with 2,000 training samples. Further, diverse generations by our method are also shown. High motion regions are indicated by a red bounding box, while spatial regions exhibiting high diversity are shown by a green bounding box. Note scaled uncertainty higher than 0.05 is shown in red.

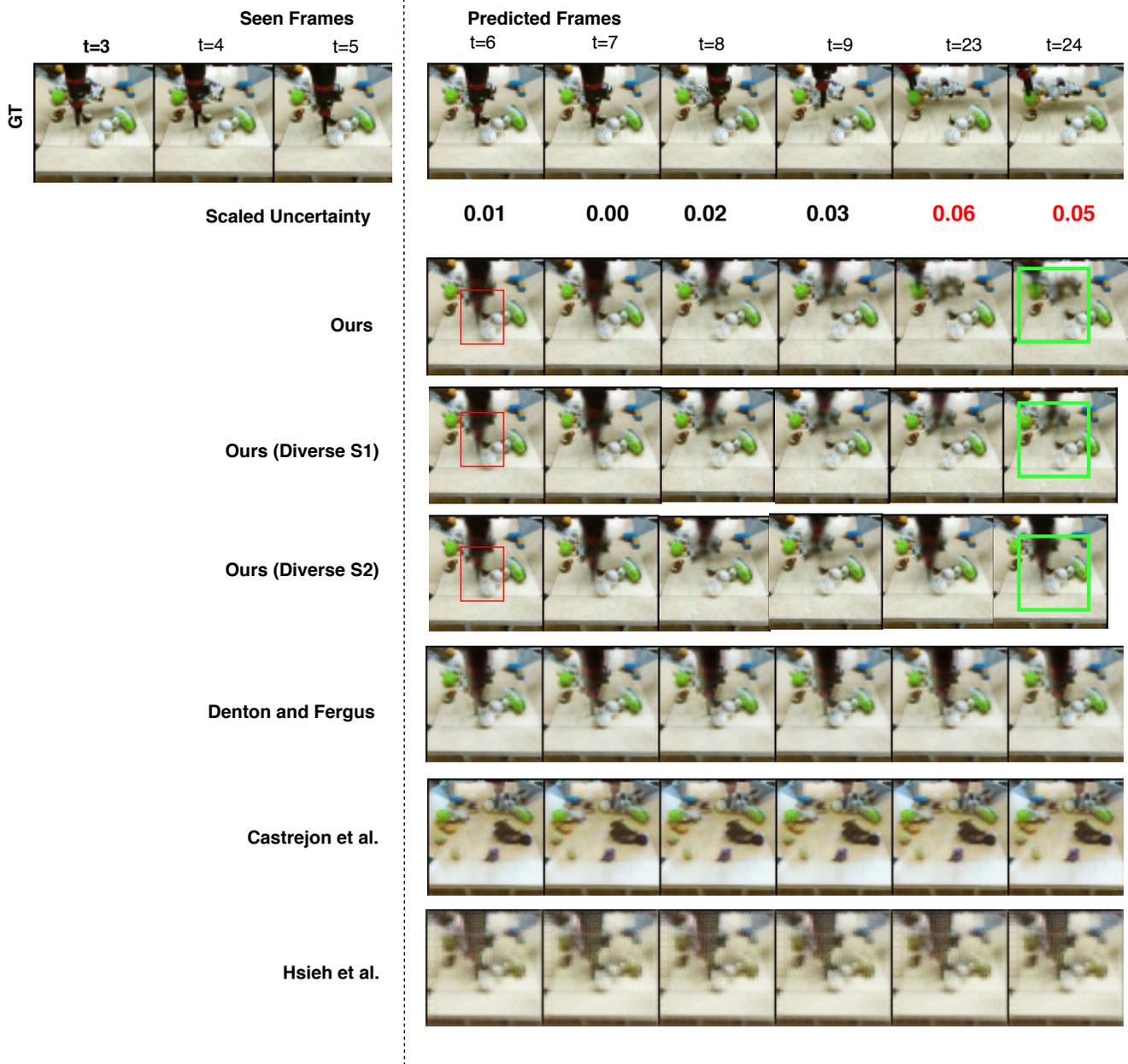


Figure 24. Visualization of generations by our method versus competing baselines on the BAIR Robot Push Dataset, trained with 2,000 training samples. Further, diverse generations by our method are also shown. High motion regions are indicated by a red bounding box, while spatial regions exhibiting high diversity are shown by a green bounding box. Note scaled uncertainty higher than 0.05 is shown in red.

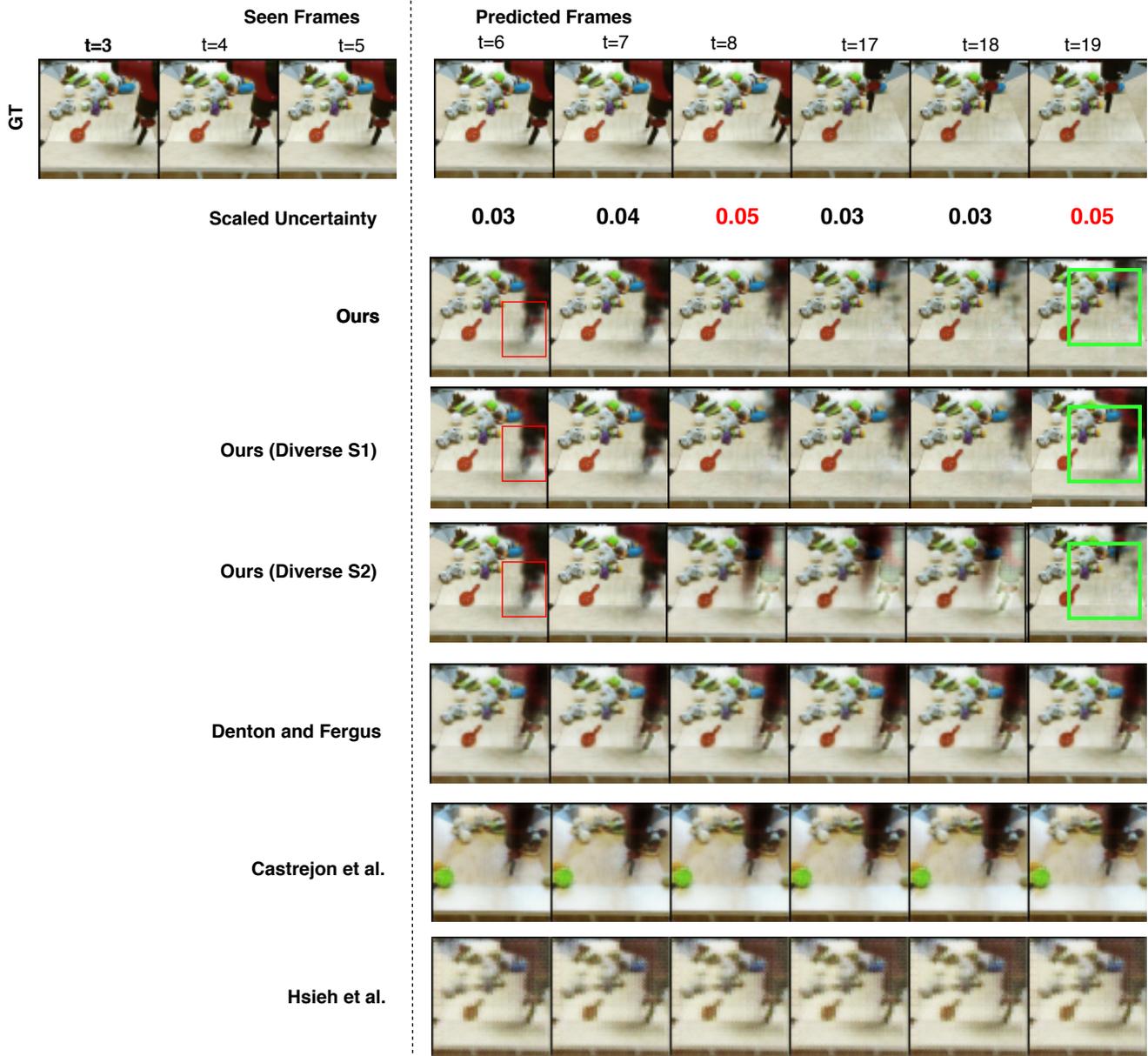


Figure 25. Visualization of generations by our method versus competing baselines on the BAIR Robot Push Dataset, trained with 2,000 training samples. Further, diverse generations by our method are also shown. High motion regions are indicated by a red bounding box, while spatial regions exhibiting high diversity are shown by a green bounding box. Note scaled uncertainty higher than 0.05 is shown in red.

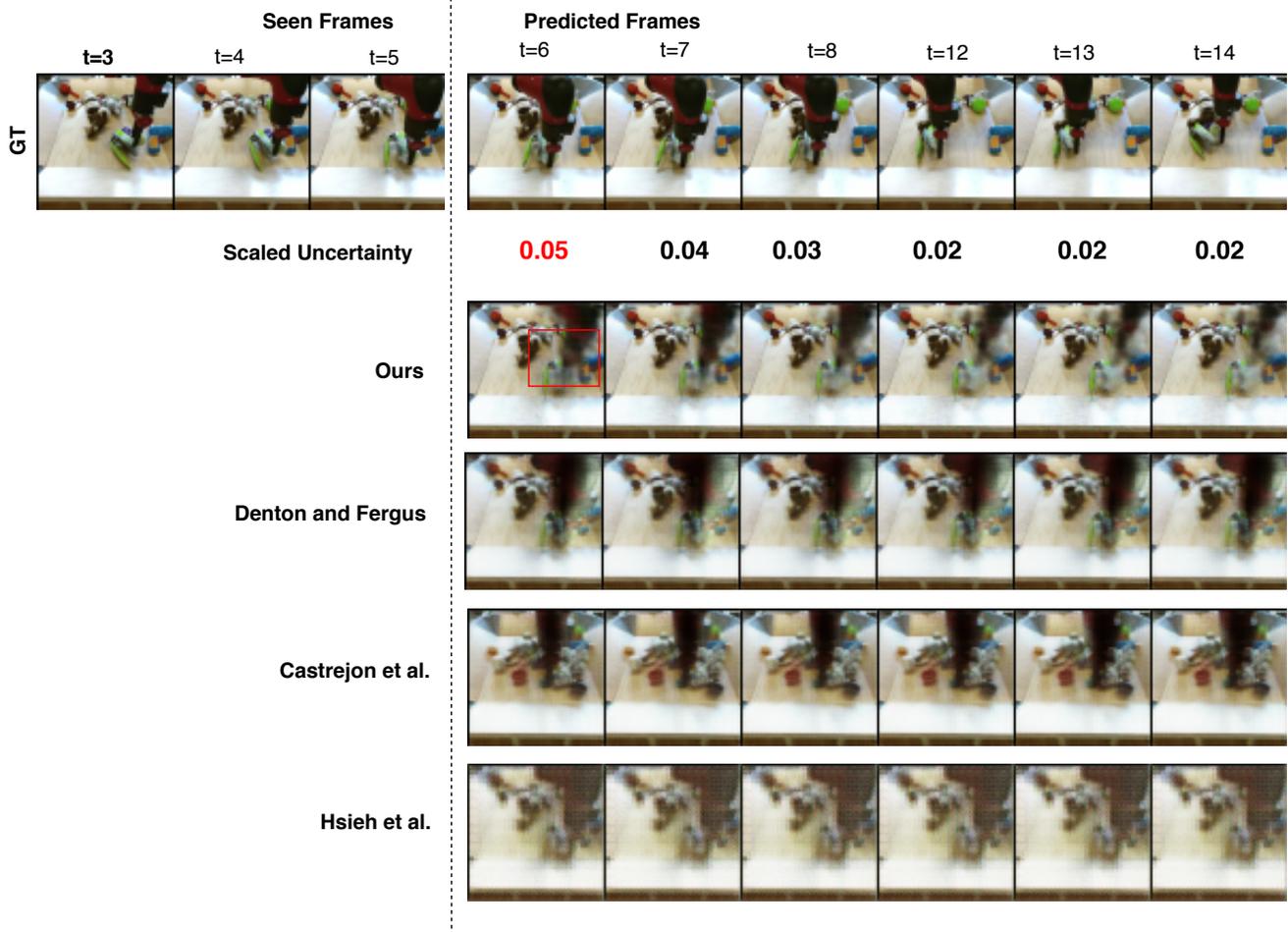


Figure 26. Visualization of generations by our method versus competing baselines on the BAIR Robot Push Dataset, trained with 2,000 training samples. High motion regions are indicated by a red bounding box. Note scaled uncertainty higher than 0.05 is shown in red.

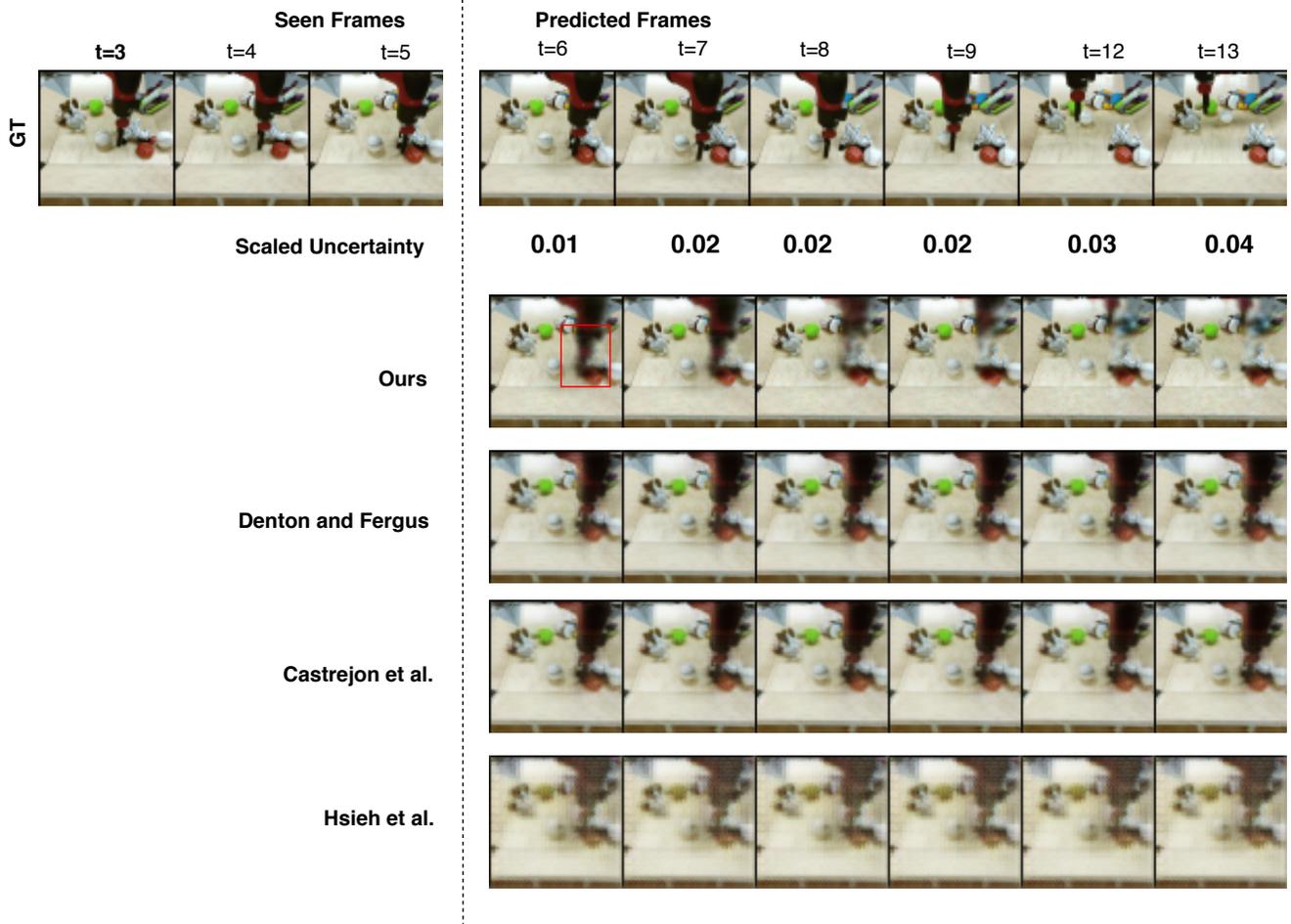


Figure 27. Visualization of generations by our method versus competing baselines on the BAIR Robot Push Dataset, trained with 2,000 training samples. High motion regions are indicated by a red bounding box. Note scaled uncertainty higher than 0.05 is shown in red.

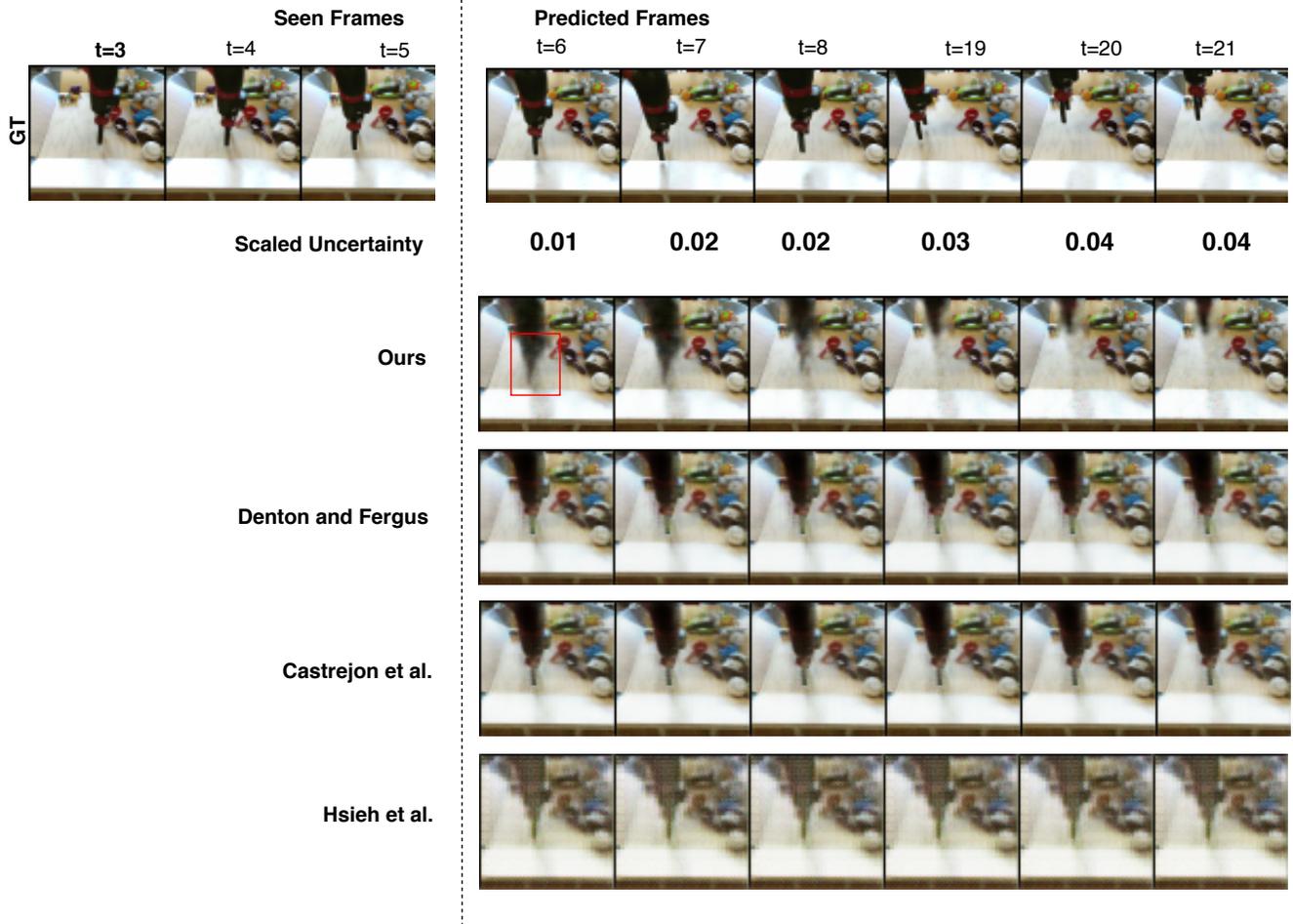


Figure 28. Visualization of generations by our method versus competing baselines on the BAIR Robot Push Dataset, trained with 2,000 training samples. High motion regions are indicated by a red bounding box. Note scaled uncertainty higher than 0.05 is shown in red.

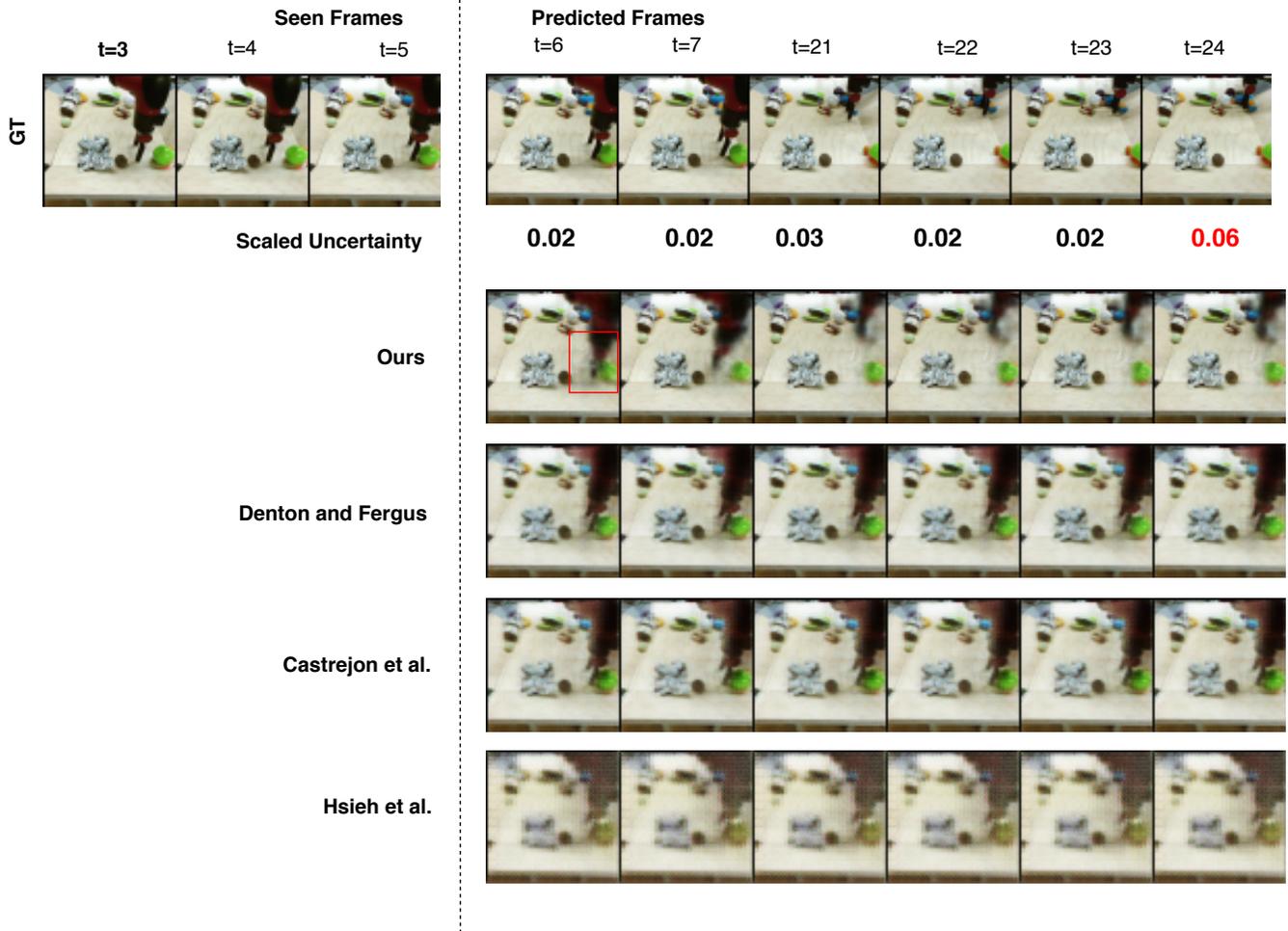


Figure 29. Visualization of generations by our method versus competing baselines on the BAIR Robot Push Dataset, trained with 2,000 training samples. High motion regions are indicated by a red bounding box. Note scaled uncertainty higher than 0.05 is shown in red.

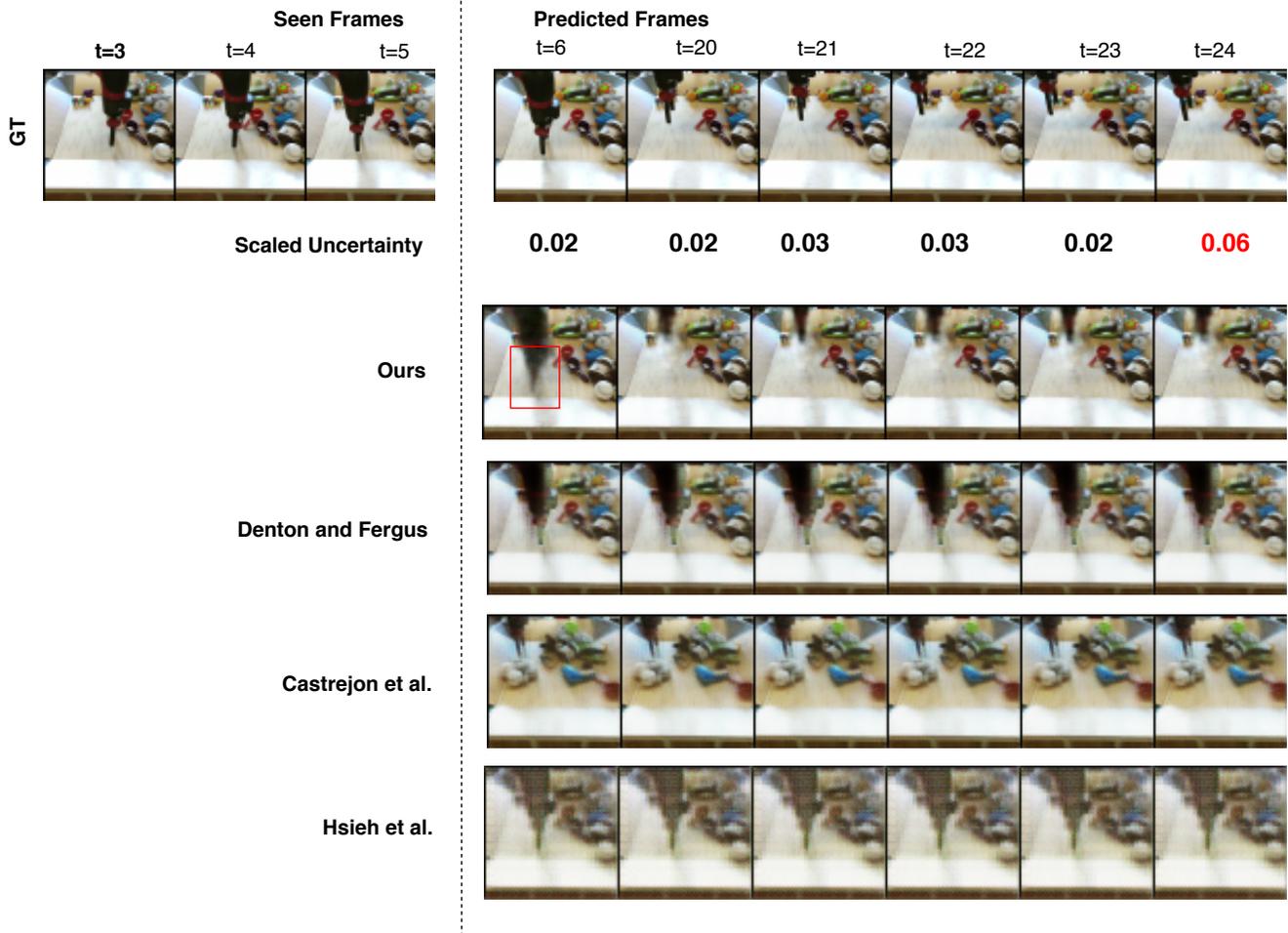


Figure 30. Visualization of generations by our method versus competing baselines on the BAIR Robot Push Dataset, trained with 2,000 training samples. High motion regions are indicated by a red bounding box. Note scaled uncertainty higher than 0.05 is shown in red.

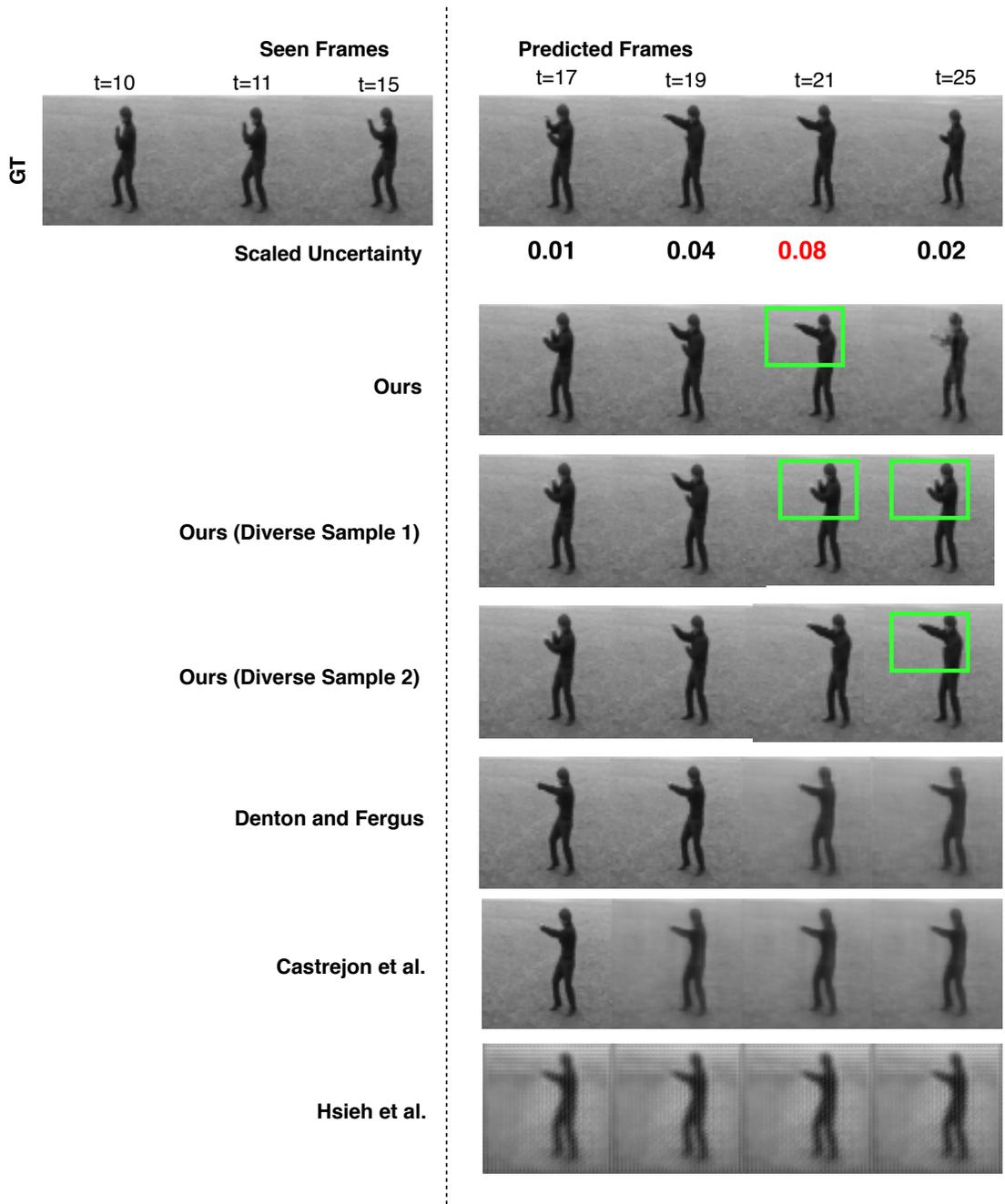


Figure 31. Visualization of generations by our method versus competing baselines on the KTH Action Dataset, trained with the full training data of 1,911 training samples. Further, diverse generations by our method are also shown. Spatial regions exhibiting high diversity are shown by a green bounding box. Note scaled uncertainty higher than 0.05 is shown in red.

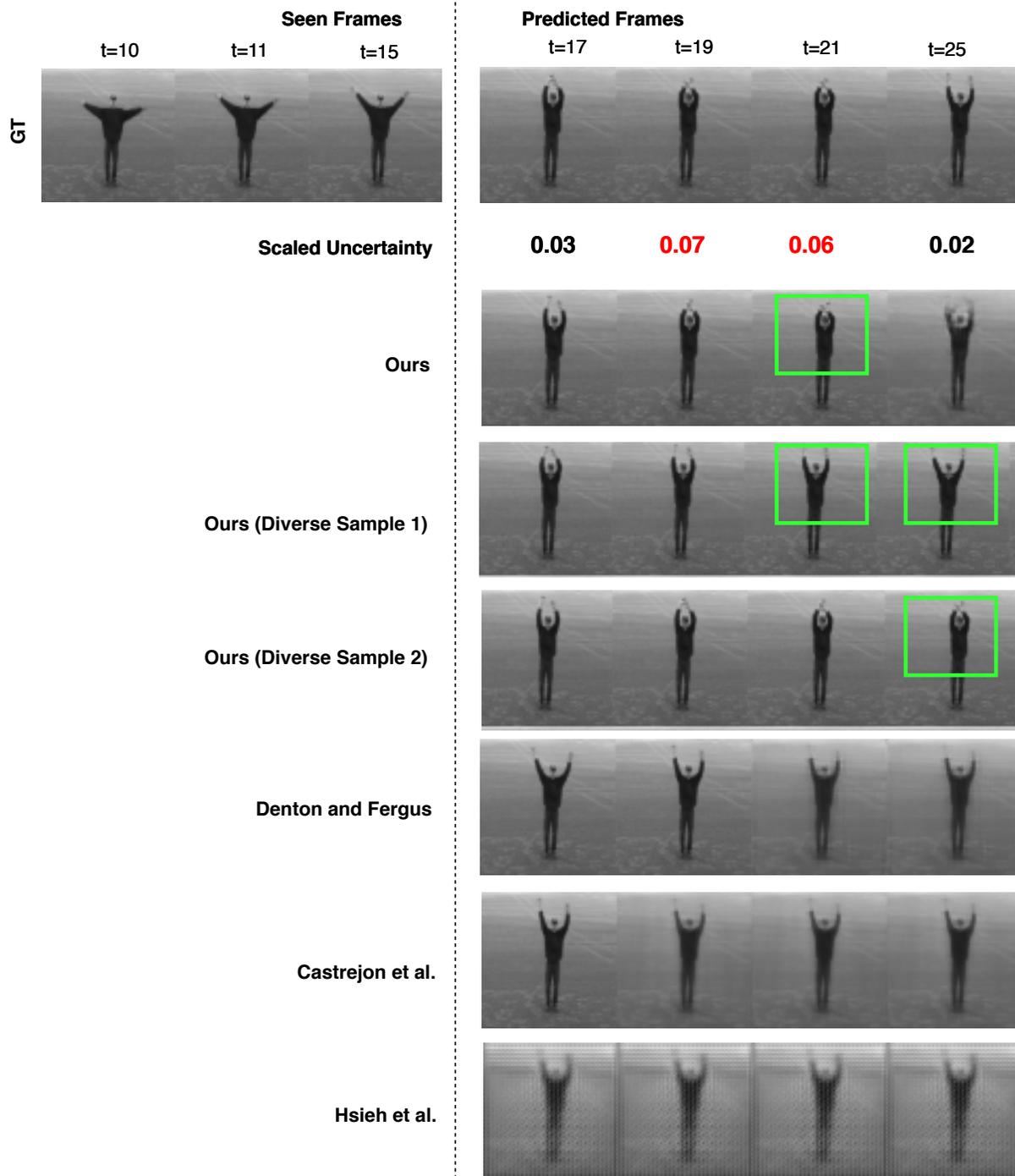


Figure 32. Visualization of generations by our method versus competing baselines on the KTH Action Dataset, trained with the full training data of 1,911 training samples. Further, diverse generations by our method are also shown. Spatial regions exhibiting high diversity are shown by a green bounding box. Note scaled uncertainty higher than 0.05 is shown in red.

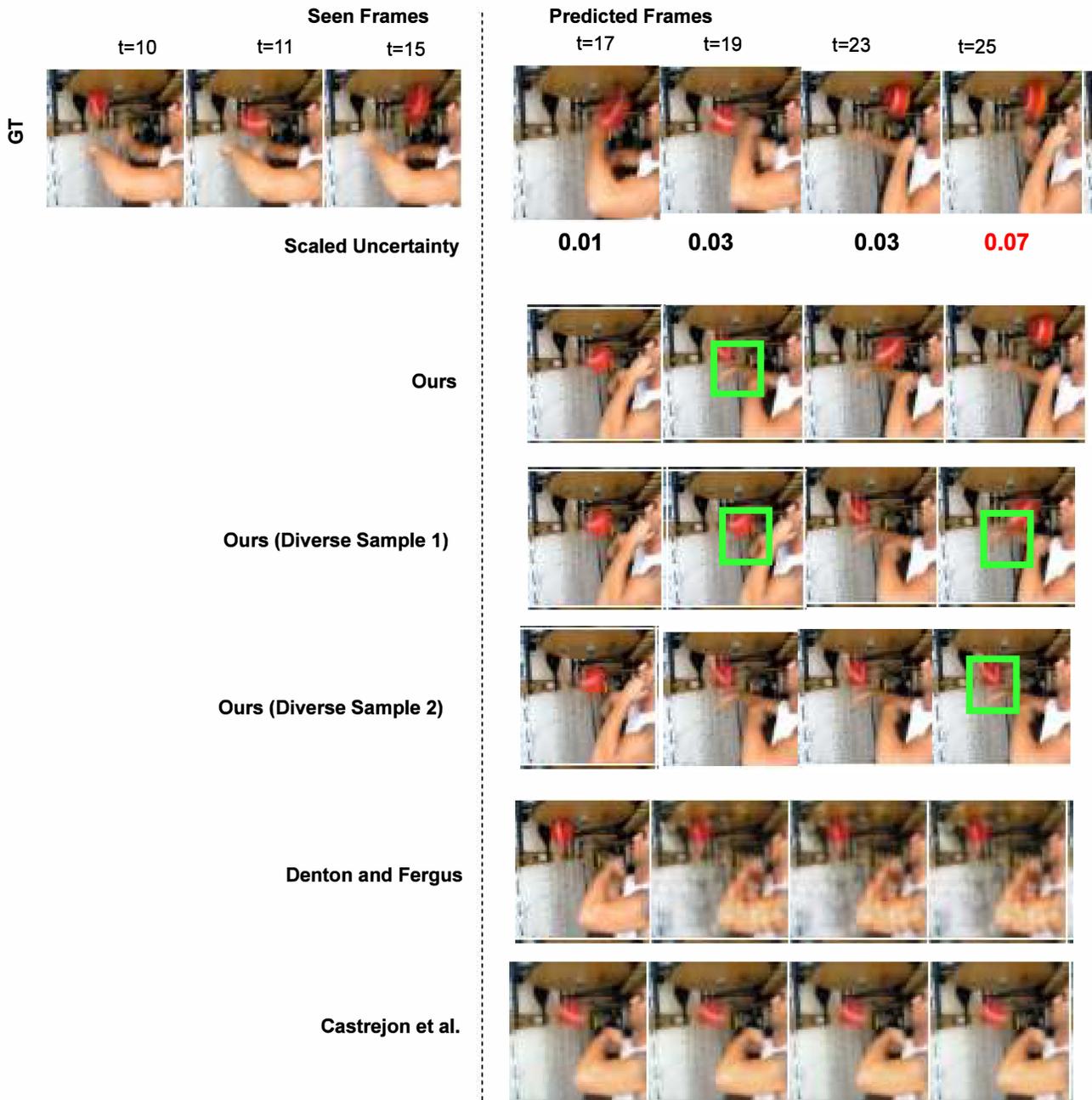


Figure 33. Visualization of generations by our method versus competing baselines on the UCF-101 Dataset, trained with the full training data of 11,425 training samples. Further, diverse generations by our method are also shown. Spatial regions exhibiting high diversity are shown by a green bounding box. Note scaled uncertainty higher than 0.05 is shown in red.

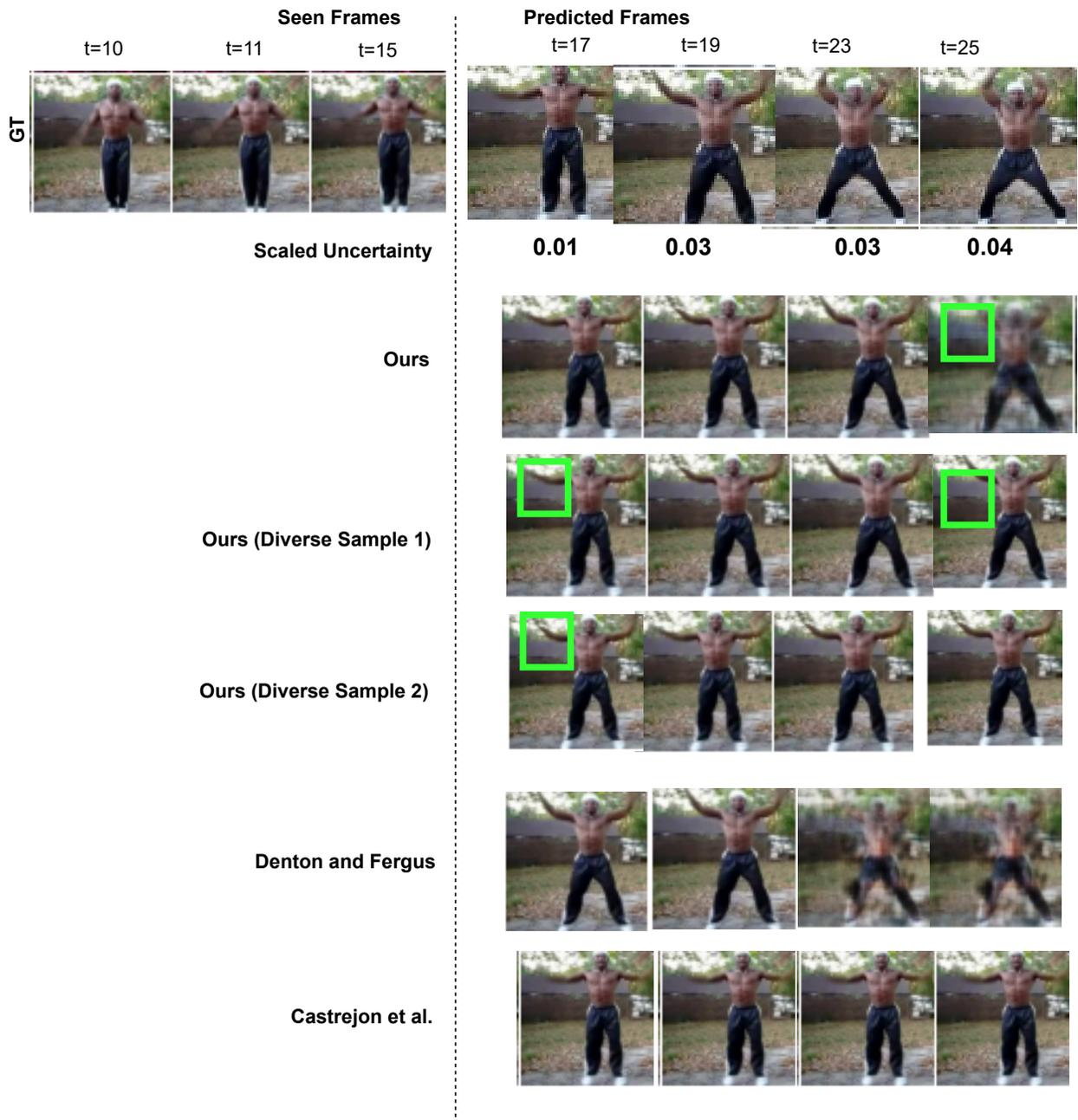


Figure 34. Visualization of generations by our method versus competing baselines on the UCF-101 Dataset, trained with the full training data of 11,425 training samples. Further, diverse generations by our method are also shown. Spatial regions exhibiting high diversity are shown by a green bounding box. Note scaled uncertainty higher than 0.05 is shown in red.