# Latent Space Subdivision:
# Stable and Controllable Time Predictions for Fluid Flow

**Steffen Wiewel** [1] **Byungsoo Kim** [2] **Vinicius C. Azevedo** [2] **Barbara Solenthaler** [2] **Nils Thuerey** [1]

## Abstract

We propose an end-to-end trained neural network architecture to robustly predict the complex dynamics of fluid flows with high temporal stability. We focus on single-phase smoke simulations in 2D and 3D based on the incompressible Navier-Stokes (NS) equations, which are relevant for a wide range of practical problems. To achieve stable predictions for long-term flow sequences, a convolutional neural network (CNN) is trained for spatial compression in combination with a temporal prediction network that consists of stacked Long Short-Term Memory (LSTM) layers. Our core contribution is a novel latent space subdivision (LSS) to separate the respective input quantities into individual parts of the encoded latent space domain. This allows to distinctively alter the encoded quantities without interfering with the remaining latent space values and hence maximizes external control. By selectively overwriting parts of the predicted latent space points, our proposed method is capable to robustly predict long-term sequences of complex physics problems. In addition, we highlight the benefits of a recurrent training on the latent space creation, which is performed by the spatial compression network.

## 1. Introduction

Computing the dynamics of fluids requires solving a set of complex equations over time. This process is computationally very expensive, especially when considering that the stability requirement poses a constraint on the maximal time step size that can be used in a simulation.

Due to the high computational resources, approaches for machine learning based physics simulations have recently

---
[1]Department of Computer Science, Technical University of Munich, Germany [2]Department of Computer Science, ETH Zurich, Switzerland. Correspondence to: Steffen Wiewel <wiewel@in.tum.de>.

been explored. One of the first approaches used Regression Forest as a regressor to forward the state of a fluid over time (Ladický et al., 2015). Handcrafted features have been used, representing the individual terms of the Navier-Stokes equations. These context-based integral features can be evaluated in constant time and robustly forward the state of the system over time. In contrast, using neural networks for the time prediction has the advantage that no features have to be defined manually, and hence these methods have recently gained increased attention. In graphics, the presented neural prediction methods (Wiewel et al., 2019; Kim et al., 2019; Morton et al., 2018) use a two-step approach, where first the physics fields are translated into a compressed representation, i.e., the latent space. Then, a second network is used to predict the state of the system over time in the latent space. The two networks are trained individually, which is an intuitive approach as spatial and temporal representations can be separated by design. In practice, the first network (i.e., the autoencoder) introduces small errors in the encoding and decoding in each time step. In combination with a temporal prediction network these errors accumulate over time, introducing drifting over prolonged time spans and can even lead to instability, as seen in Figure 4 (top right). This is especially problematic in supervised learned latent space representations, since the drift will shift the initial, user-specified conditions (e.g., an object's position) into an erroneous latent space configuration originated from different conditions.

Like previous work, we use a neural network to predict the motion of a fluid over time, but with the central goal to increase accuracy and robustness of long-term predictions. We propose to use a joint end-to-end training of both components, the fluid state compression and temporal prediction of the motion. A key observation is also the need to control the learned latent space, such that we can influence it to impose boundary conditions and other known information that is external to the simulation state. We therefore propose a latent space subdivision that separates the encoded quantities in the latent space domain. The subdivision is enforced with a latent space split soft-constraint for the input quantities velocity and density, and allows to alter the individual encoded components separately. This separation is a key component to robustly predict long-term sequences.

## 2. Related Work and Background

Our work concentrates on single-phase flows, which are usually modeled by a pressure-velocity formulation of the incompressible Navier-Stokes equations:

$$\frac{\partial \boldsymbol{u}}{\partial t} + \boldsymbol{u} \cdot \nabla \boldsymbol{u} = \frac{-1}{\rho} \nabla p + \nu \nabla^2 \boldsymbol{u} + \boldsymbol{g} \qquad (1)$$

$$\nabla \cdot \boldsymbol{u} = 0,$$

where $p$ denotes pressure, $\boldsymbol{u}$ is the flow velocity, and $\rho, \nu, \boldsymbol{g}$ denote density, kinematic viscosity, and external forces, respectively. The density $\rho$ is passively advected by the velocity $\boldsymbol{u}$. A complete overview of fluid simulation techniques in computer graphics can be found in (Bridson, 2015).

Data-driven flow modelling encompasses two distinguished and complementary efforts: dimensionality reduction and reduced-order modelling. Dimensionality reduction (e.g., Singular Value Decomposition) scales down the analyzed data into a set of important features in an attempt to increase the sparsity of the representation, while reduced-order modelling (e.g., Galerkin Projection) describes the spatial and temporal dynamics of a system represented by a set of reduced parameters. In computer graphics, the work of Treuille et al. (2006) was the first to use Principal Component Analysis (PCA) for dimensionality reduction coupled with a Galerkin Projection method for subspace simulation. This approach was later extended (Kim & Delaney, 2013) with a cubature approach for enabling Semi-Lagrangian and Mac-Cormack (Selle et al., 2008) advection schemes, while improving the handling of boundary conditions. Reduced-order modelling was also explored to accelerate the pressure projection step in liquid simulations (Ando et al., 2015).

Instead of computing reduced representations from pre-simulated velocity fields, alternative basis functions can be used for reduced-order modelling; examples of basis functions include Legendre Polynomials (Gupta & Narasimhan, 2007), modular (Wicke et al., 2009) and spectral (Long & Reinhard, 2009) representations. Also Laplacian Eigenfunctions have been successfully employed for dimensionality reduction, due to their natural sparsity and inherent incompressibility. De Witt et al. (2012) combined Laplacian Eigenfunctions with a Galerkin Projection method, enabling fast and energy-preserving fluid simulations. The approach was extended to handle arbitrarily-shaped domains (Liu et al., 2015), combined with a Discrete Cosine Transform (DCT) for compression (Jones et al., 2016), and improved for scalability (Cui et al., 2018).

The aforementioned methods for data-driven flow modelling use linear basis functions for dimensionality reduction. This enables the use of Galerkin Projection for subspace integration, but it limits the power of the reconstruction when compared to non-linear embeddings. The latent spaces generated by autoencoders (AE) are non-linear and richly cap-

ture the input space with fewer variables (Radford et al., 2016; Wu et al., 2016). In light of that, Wiewel et al. (2019) combined a latent space representation with recurrent neural networks (RNN) to predict the temporal evolution of fluid functions in the latent space domain. Kim et al. (2019) introduced a generative deep neural network for parameterized fluid simulations that only takes a small set of physical parameters as input to very efficiently synthesize points in the learned parameter space. Their method also proposes an extension to latent space integration by training a fully connected neural network that maps subsequent latent spaces. Our work is related to these two methods, but a main difference is that we use an end-to-end training of both the spatial compression and the temporal prediction. In combination with our latent space subdivision, our predictions are more stable, while previous approaches fail to properly recover long-term integration correspondences due to the lack of autoencoder regularization.

For particle-based fluid simulations, a temporal state prediction using Regression Forest was presented in Ladicky et al. (2015). Handcrafted features are evaluated in particle neighborhoods and serve as input to the regressor, which then predicts the particle velocity of the next time step. Machine learning has also been used in the context of grid-based (Eulerian) fluid simulations. Tompson et al. (2017) used a convolutional neural network (CNN) to model spatial dependencies in conjunction with an unsupervised loss function formulation to infer pressure fields. A simpler three-layer fully connected neural network for the same goal was likewise proposed (Yang et al., 2016). As an alternative, learned time evolutions for Koopman operators were proposed (Morton et al., 2018), which however employ a pre-computed dimensionality reduction via PCA. Chu et al. (2017) enhance coarse fluid simulations to generate highly detailed turbulent flows. Individual low-resolution fluid patches were tracked and mapped to high-resolution counterparts via learned descriptors. Xie et al. (2018) extended this approach by using a conditional generative adversarial network with a spatio-temporal discriminator supervision. Small-scale splash details in hybrid fluid solvers were targeted with deep learning-based stochastic models (Um et al., 2018). For a review of machine learning applied to fluid mechanics we refer the reader to (Brunton et al., 2019).

In the context of data-driven flow modelling, methods to interpolate between existing data have been presented. Raveendran et al. (2014) presented a technique to smoothly blend between pre-computed liquid simulations, which was later extended with more controllability (Manteaux et al., 2016). Thuerey et al. (2016) used dense space-time deformations represented by grid-based signed-distance functions for interpolation of smoke and liquids, and Sato et al. (2018) interpolated velocity fields by minimizing an energy functional.

# 3. Method

The central goal of our models is to robustly and accurately predict long-term sequences of flow dynamics. For this, we need an autoencoder to translate high-dimensional physics fields into a compressed representation (latent space) and a temporal prediction network to advance the state of the simulation over time. A key observation is that if these two network components are trained individually, neither component has a holistic view on the underlying problem. The autoencoder, consisting of an encoder $E$ and a decoder $D$, generates a compressed representation $\boldsymbol{c} = E(\boldsymbol{x})$, which focuses solely on the reconstruction $\tilde{\boldsymbol{x}} = D(\boldsymbol{c})$ of the given input $\boldsymbol{x}$. Hence, the loss function to minimize is given by $\|\boldsymbol{x} - \tilde{\boldsymbol{x}}\|$. Without considering the aspect of time, the autoencoder's latent space only stores spatial descriptors. Due to the exclusive focus on space, temporally consecutive data points are not necessarily placed close to each other in the latent space domain. This poses substantial challenges for the temporal prediction network.

Therefore, we consider the aspect of time within the training of the autoencoder in order to shape its latent space with respect to temporal information, in addition to the spatial information. Thus, we propose an end-to-end training procedure, where we train our autoencoder and temporal prediction network simultaneously by internally connecting the latter as a recurrent block to the encoding and decoding blocks of the spatial autoencoder. As a result, the latent space domain is aware of temporal changes, and can yield temporally coherent latent space points that are suitable for the time prediction network. By default, our training process includes the combined training of our spatial autoencoder and temporal prediction network as shown in Figure 1. In this figure, the encoder E, decoder D and prediction network P are duplicated for visualization purposes. In the next sections, we describe each individual network in more detail.

## 3.1. Spatial Encoding

The spatial encoding of the data is performed by a regular autoencoder, the network for which is split into an encoding and decoding part (Kim et al., 2019). The encoder contains 16 convolution layers with skip connections, which connect its internal layers, followed by one fully-connected layer. The decoder consists of a fully-connected layer, which is followed by 17 convolution layers with skip connections. For the fluid simulation dataset used in this work, the input is either 2D or 3D, leading to the usage of 2D- or 3D-convolutions and a feature dimension of 3 or 4, respectively. Furthermore, a data-specific curl-layer is appended to the decoder network to enforce zero divergence in the resulting velocity field (Kim et al., 2019), as required by the NS equations (see Equation 1).
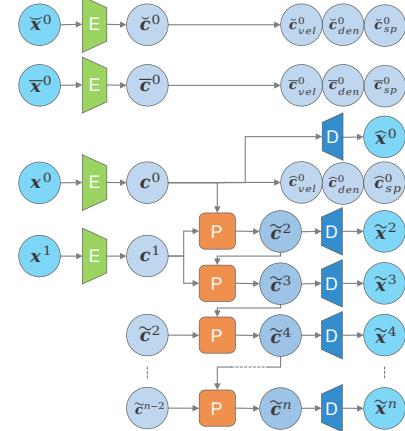


Figure 1: Combined training of autoencoder and temporal prediction. The prediction networks input window is set to $w = 2$. Thus, the count of recurrent predictions is $n_i = n - 2$. The LSS is enforced by applying the split loss on $\bar{\boldsymbol{c}}_{den}^0$ and $\breve{\boldsymbol{c}}_{vel}^0$ for velocity and density, respectively. A direct AE reconstruction loss is only performed on $\hat{\boldsymbol{x}}^0$.

The dimensionality of the latent space $\boldsymbol{c}$ for a given $\boldsymbol{x}$ is defined by the final layer of the encoder and can be freely chosen. We pass a velocity $\boldsymbol{u}$ as well as a density field $\rho$ to our encoder, i.e., $\boldsymbol{x} = [\boldsymbol{u}, \rho]$. The velocity field is an active quantity that is used in fluid simulations to advect a passive quantity forward in time. In case of smoke simulations, which is a specific instance of fluid simulations, the passive density field is advected by the flow velocity. As a result, $\boldsymbol{c}$ contains information about both the active and passive fields, i.e., velocity and density. Hence, we can accurately advect the passive quantity density with a velocity field with low computational effort, it makes sense to compute the advection outside of the network, and project the new state into the latent space.

In order to be able to alter active and passive fields individually in the compressed representation $\boldsymbol{c}$, with given input field $\boldsymbol{x} = [\boldsymbol{x}_{vel}, \boldsymbol{x}_{den}]$ where the subscripts $vel$ and $den$ thereby denote the velocity and density part, we subdivide $\boldsymbol{c}$ into separate parts for velocity $\boldsymbol{c}_{vel}$ and density $\boldsymbol{c}_{den}$, respectively. This property of the latent space is needed for projecting the new state of the passive quantity into the latent space domain. Additionally, to exert explicit external control over the prediction, we designate another part of $\boldsymbol{c}$ to contain supervised parameters, called $\boldsymbol{c}_{sp}$ (Kim et al., 2019). In our case of, e.g., a smoke simulation with a rotating cup filled with smoke, such supervised parameters can be the position of a smoke source or the rotation angle of a solid obstacle. With this subdivision, we increase stability of our predictions and allow for explicit external control. This subdivision is described in Equation 2, where $v$, $d$, and $sp$ describe the indices of the velocity, density, and supervised

parameter parts in the latent space domain, respectively.

$$\boldsymbol{c} = \begin{bmatrix} \boldsymbol{c}_{vel} \,\big\backslash\, \boldsymbol{c}_{den} \,\big|\, \boldsymbol{c}_{sp} \end{bmatrix} \qquad (2)$$
$$\text{where} \quad \boldsymbol{c}_{vel} = \begin{bmatrix} c_0, \dots, c_v \end{bmatrix},$$
$$\boldsymbol{c}_{den} = \begin{bmatrix} c_{v+1}, \dots, c_d \end{bmatrix},$$
$$\boldsymbol{c}_{sp} = \begin{bmatrix} c_{d+1}, \dots, c_{sp} \end{bmatrix}.$$

To arrive at the desired subdivision, the split loss $\mathcal{L}_{split}$ (see Equation 3) is used as a loss function in the training process. It is modelled as a soft constraint and thereby does not enforce the parts $\boldsymbol{c}_{vel}$ and $\boldsymbol{c}_{den}$ to be strictly disjunct. The loss is defined as

$$\mathcal{L}_{split}(\boldsymbol{c}, I_s, I_e) = \sum_{i=I_s}^{I_e} \|c_i\|_1. \qquad (3)$$

Since we divide the latent space in three parts, $\mathcal{L}_{split}$ is applied twice. For the velocity part $\boldsymbol{c}_{vel}$, the indices $I_s = v+1$ and $I_e = d$ are chosen to indicate that the density part must not be used on encoding velocities, i.e., $\mathcal{L}_{split}(\boldsymbol{c}, v+1, d)$. In contrast to the previous limits, for the density part $\boldsymbol{c}_{den}$ the velocity part is indicated by choosing the indices $I_s = 0$ and $I_e = v$, i.e. $\mathcal{L}_{split}(\boldsymbol{c}, 0, v)$. First, only the velocity part $\boldsymbol{x}_{vel}$ of input $\boldsymbol{x}$ is encoded (the density part $\boldsymbol{x}_{den}$ is zero, i.e., $\bar{\boldsymbol{x}} = [\boldsymbol{x}_{vel}, 0]$), yielding $\bar{\boldsymbol{c}}_{vel}$. Vice versa, the density part $\boldsymbol{x}_{den}$ of input $\boldsymbol{x}$ is encoded, whereas the velocity part $\boldsymbol{x}_{vel}$ is replaced with zeros, i.e., $\check{\boldsymbol{x}} = [0, \boldsymbol{x}_{den}]$, resulting in $\check{\boldsymbol{c}}_{den}$. Therefore, the loss is applied twice for the $\bar{\boldsymbol{c}}$ and $\check{\boldsymbol{c}}$ encodings as $\mathcal{L}_{split}(\bar{\boldsymbol{c}}, v+1, d)$ and $\mathcal{L}_{split}(\check{\boldsymbol{c}}, 0, v)$, respectively.

In order to exhibit external control over the prediction, $\boldsymbol{c}_{sp}$ is enforced to contain parameters describing certain attributes of the simulation. While training the network, an additional soft-constraint is applied, which forces the encoder to produce the supervised parameters. The soft-constraint is implemented as the mean-squared error of the values generated by the encoder $\hat{\boldsymbol{c}}_{sp}$ and the ground truth data $\boldsymbol{c}_{sp}$ and consitutes the supervised loss $\mathcal{L}_{sup}$ as

$$\mathcal{L}_{sup}(\boldsymbol{c}_{sp}, \hat{\boldsymbol{c}}_{sp}) = \|\boldsymbol{c}_{sp} - \hat{\boldsymbol{c}}_{sp}\|_2^2. \qquad (4)$$

Additionally, an AE loss $\mathcal{L}_{AE}$ (Equation 5) is applied to the decoded field $\hat{\boldsymbol{x}}$. It forces the velocity part of the decoded field to be close to the input velocity by applying the mean-absolute error. To take the rate of change of the velocities into consideration as well, the mean-absolute error of the velocities' gradient is added to the formulation. In contrast, the density part is handled by directly applying the mean-squared error on the decoded output density and the input. The AE loss is thereby defined as

$$\mathcal{L}_{AE}(\boldsymbol{x}, \hat{\boldsymbol{x}}) = \|\boldsymbol{x}_{vel} - \hat{\boldsymbol{x}}_{vel}\|_1$$
$$+ \|\nabla \boldsymbol{x}_{vel} - \nabla \hat{\boldsymbol{x}}_{vel}\|_1$$
$$+ \|\boldsymbol{x}_{den} - \hat{\boldsymbol{x}}_{den}\|_2^2. \qquad (5)$$

## 3.2. Time Prediction Network

The prediction network performs a temporal transformation of its input to the temporal consecutive state. The inputs are a series of $w$ consecutive input states. The prediction network block contains two recurrent LSTM layers, followed by two 1D-convolutional layers. In our case of 2D smoke simulations, two consecutive latent space points of dimension 16 are used as input. Those are fed to the prediction layers and result in one latent space point of dimension 16, called the residuum $\Delta \boldsymbol{c}^t$. Afterwards, the residuum is added to the last input state to arrive at the next consecutive state, i.e., $\tilde{\boldsymbol{c}}^{t+1} = \boldsymbol{c}^t + \Delta \boldsymbol{c}^t$.

Due to the subdivision capability of our autoencoder, our temporal prediction network supports external influence over the predictions it generates. After each prediction, it is possible to replace or update information without the need of re-encoding the predicted data. Instead, only parts of the predicted latent space point can be replaced, enabling fine-grained control over the flow. For example, in the case of smoke simulations, the passive smoke density quantity can be overwritten with an externally updated version, i.e., the $\boldsymbol{c}_{den}$ part is replaced by $\boldsymbol{c}$. This allows for adding new smoke sources or modifying the current flow by removing smoke from certain parts of the simulation domain.

Considering the exposition of the prediction input window $w$, which can be chosen freely, and the desired internal iteration count $n_i = n - w$, the additive prediction error is brought into consideration for the prediction network P while training, i.e., it is traversed $n_i$ times. This leads to a combined training loss of AE and P defined as

$$\mathcal{L} = \mathcal{L}_{AE,direct} + \mathcal{L}_{sup} + \mathcal{L}_{split,vel} + \mathcal{L}_{split,den}$$
$$+ \sum_{n=0}^{n_i} \big( \mathcal{L}_{AE,pred_{n_i}} \big), \qquad (6)$$

where $\mathcal{L}_{AE}$ is applied to the corresponding pairs of the decoded outputs $\tilde{\boldsymbol{x}}^2 ... \tilde{\boldsymbol{x}}^n$ of the prediction network P and their corresponding ground-truth $\boldsymbol{x}^2 ... \boldsymbol{x}^n$. Thereby, our final loss is the sum of all the previously presented losses.

In our combined training approach, both networks update their weights by applying holistic gradient evaluations, i.e., are trained end-to-end. The benefit of the end-to-end training is that the spatial autoencoder AE also incorporates temporal aspects when updating its weights. In addition, by recurrently applying $\mathcal{L}_{AE}$ on the predictions, the prediction network P is trained to actively minimize accumulating additive errors. To incorporate temporal awareness in the autoencoder, the decoder block is connected to the individual prediction outputs and is thereby reused several times in one network traversal. The recurrent usage of the decoding block is commonly known as weight sharing (Bromley et al., 1993). Furthermore, by applying the prediction losses

on the decoded predictions, the spatial autoencoder adapts to the changes induced by the temporal prediction as well, which furthers the focus of the autoencoder to produce latent spaces suitable for temporal predictions. As a result, the prediction network is capable of robustly and accurately predicting long-term sequences of complex fluid flows.

## 4. Training Datasets

The datasets we used to train our networks contain randomized smoke flows simulated with an open source framework (Thuerey & Pfaff, 2018). In total, three different scene setups were used to capture a wide range of complex physical behavior. The first scene contains a moving smoke inflow source that generates hot smoke continuously, which is rising and producing complex swirls (see Figure 2).
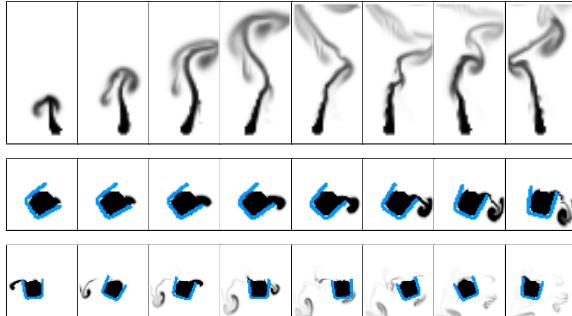


Figure 2: Example sequences of our 2D datasets: moving smoke (top), rotating (center) and moving cup (bottom). The smoke density is shown as black and the cup-shaped obstacle in blue.

The second and third scenes simulate cold smoke in a cup-shaped obstacle. The former rotates the cup randomly around a fixed axis, while the latter additionally applies a translation (see Figure 2). The rising smoke and the rotating cup scene each expose one control parameter, i.e., movement on the x-axis and rotation around the z-axis, whereas the rotating and moving cup scene exposes both of these control parameters. Each of the three datasets in 2D contains 200 randomly generated scenes with 600 consecutive frames. Additionally, the moving smoke as well as the rotating and moving cup dataset was generated in 3D with 100 randomly generated scenes and 600 consecutive frames (see Table 1).

Table 1: Statistics of our datasets.

| Scene Type | Resolution | # Scene | # Frames |
|---|---|---|---|
| Rotating and Moving Cup (3D) | $48^3$ | 100 | 600 |
| Moving Smoke (3D) | $48^3$ | 100 | 600 |
| Rotating and Moving Cup (2D) | $64^2$ | 200 | 600 |
| Rotating Cup (2D) | $48^2$ | 200 | 600 |
| Moving Smoke (2D) | $32x64$ | 200 | 600 |

## 5. Evaluation

In this section, we compare our architecture to the baseline of previous work. We also perform an ablation study on different settings of our proposed architecture to compare their respective influence on the output. We compute the mean peak signal-to-noise ratio (PSNR) for all our comparisons, i.e., larger values are better. For each case, we measure accuracy of our prediction w.r.t. density and velocity in terms of PSNR for ten simulations setups that were not seen during training.

For a thorough evaluation, we supply two prediction approaches. First, we evaluate a regular prediction approach with no reinjection of physical information (denoted *VelDen*) that is in sync with previous work (Kim et al., 2019; Wiewel et al., 2019). This approach is formulated as

$$\tilde{c}^t = P(\tilde{c}^{t-2}, \tilde{c}^{t-1}), \qquad (7)$$
$$\tilde{x}^t = [\tilde{x}_{vel}^t, \tilde{x}_{den}^t] = D(\tilde{c}^t),$$

where the previously predicted latent space points $\tilde{c}^{t-2}$ and $\tilde{c}^{t-1}$ are used to evaluate the next time step $\tilde{c}^t$. Afterwards, $\tilde{c}^t$ is decoded $D(\tilde{c}^t)$ and the density part $\tilde{x}_{den}^t$ is directly displayed, i.e., no external physical information about the system state is influencing the output of our *VelDen* benchmarks.

In the second approach, we make use of our LSS to reinject the advected density into the prediction to benefit from well understood physical computations that keep our predictions stable and can be performed in a fast manner. The prediction process utilizing our LSS is denoted as *Vel* and is given as

$$\tilde{c}^t = P(\hat{c}^{t-2}, \hat{c}^{t-1}), \qquad (8)$$
$$\tilde{x}^t = [\tilde{x}_{vel}^t, \tilde{x}_{den}^t] = D(\tilde{c}^t),$$
$$\dot{x}_{den}^t = Adv(\dot{x}_{den}^{t-1}, \tilde{x}_{vel}^t),$$
$$\dot{c}^t = E([\tilde{x}_{vel}^t, \dot{x}_{den}^t]),$$
$$\hat{c}^t = [\tilde{c}_{vel}^t, \dot{c}_{den}^t],$$
$$\hat{x}^t = [\tilde{x}_{vel}^t, \dot{x}_{den}^t],$$

where we are using the decoded predicted velocity $\tilde{x}_{vel}^t$ to advect the simulation density $\dot{x}_{den}^{t-1}$ and reinject its encoded form into our latent space $\hat{c}^t$. The new latent space representation $\hat{c}^t$ is thereby formed by concatenating the new encoded density $\dot{c}_{den}^t$ and the predicted encoded velocity field $\tilde{c}_{vel}^t$. By reinjecting the advected density field $\dot{x}_{den}^t$, we inform the prediction network about boundary conditions as well as other known physical information that is external to the prediction state. In the following we will ablate on different aspects of our method to evaluate their respective influence on the final results.

**Latent Space Temporal Awareness** The temporal awareness of our spatial autoencoder is evaluated in this section,

(a) No temporal constraints or super-
vised parameters.

(b) No temporal constraints but with
supervised parameters.

(c) With temporal constraints and su-
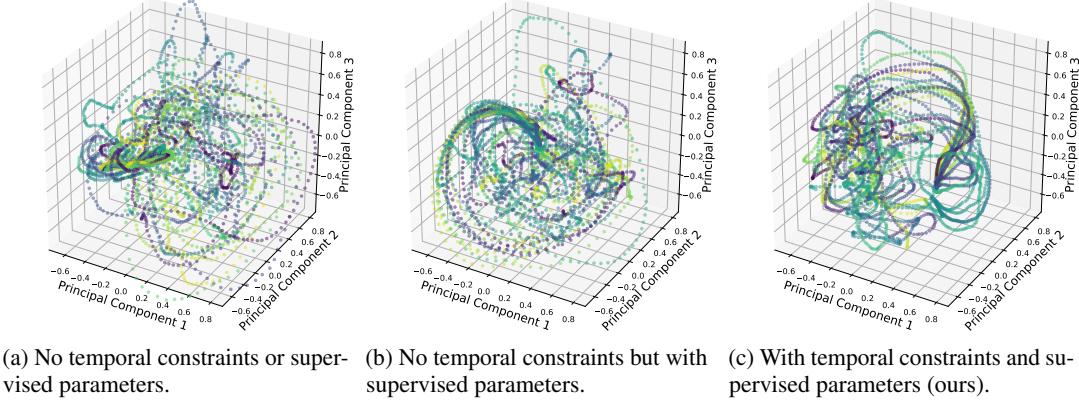pervised parameters (ours).

Figure 3: Spatial encodings of 200 frames of 20 different smoke simulations. The latent space points are normalized to
their respective maximum and processed with PCA for visualization purposes. Each color stands for a single simulation
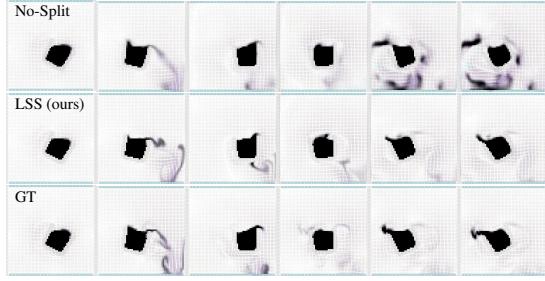and represents a series of 200 frames.



Figure 4: Long-term prediction of 400 time steps: Our
method robustly predicts the fluid flow, whereas the
regular prediction method (no-split) fails to capture
the movement and even produces unphysical behavior.
Compare the PSNR values in Table 3.



Figure 5: A circular solid obstacle, unseen during train-
ing, is placed in the upper right. Our proposed method
(center) predicts stable and realistic. The prediction
only approach (top) is not able to capture the obstacle.

since it has a significant impact on the performance of our
temporal prediction network. In Figure 3 we evaluate three
networks trained with different loss functions in terms of
the stability of the latent space they generate for sequences.
For each of the plots, 200 frames of 20 different smoke sim-
ulations were encoded to the latent space domain with an
autoencoder. The resulting latent space points were normal-
ized with their respective maximum to the range of $[-1, 1]$
and afterwards transformed to 3 dimensions with PCA. The
supervised part was removed before normalization. For our
comparison we chose 3 autoencoders with a latent space
dimension of 16.

The results in Figure 3a were generated with a classic AE
that was trained to only reproduce its input, i.e., only a
direct loss on the output (Equation 5) was applied. For
this classic AE no temporal constraints were imposed, and
no supervised parameters were added to the latent space.
The resulting PCA decomposition shows a very uneven
distribution: large distances between consecutive points
exist in some siutations, whereas a large part of the samples
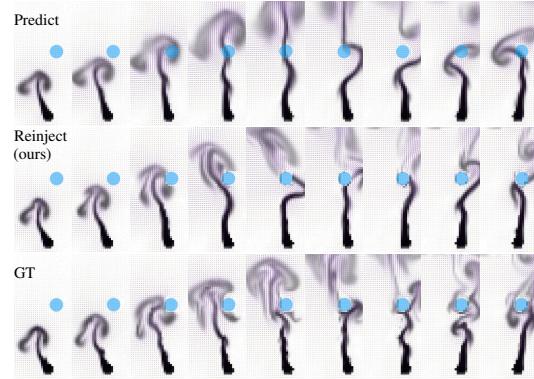are placed closely together in a cluster.

When adding the supervised parameter loss (Equation 4) in
the second evaluation (Kim et al., 2019) the trajectories be-
come more ordered, as shown in Figure 3b, but still exhibit a
noticeably uneven distribution. Thus, the supervised param-
eter, despite being excluded from the PCA, has a noticeable
influence on the latent space construction.

In Figure 3c, the results of the AE trained with our proposed
time-aware end-to-end training method are shown. This
AE applies the supervised parameter loss (Equation 4) as
well as the direct loss on the output (Equation 5) and was
trained in combination with the temporal prediction net-
work P as described in Section 4. The visualization of the
PCA decomposition shows that a strong temporal coherence
of the consecutive data points emerges. This visualization
indicates why our prediction networks yield substantial im-
provements in terms of accuracy: the end-to-end training
provides the autoencoder with gradients that guide it to-
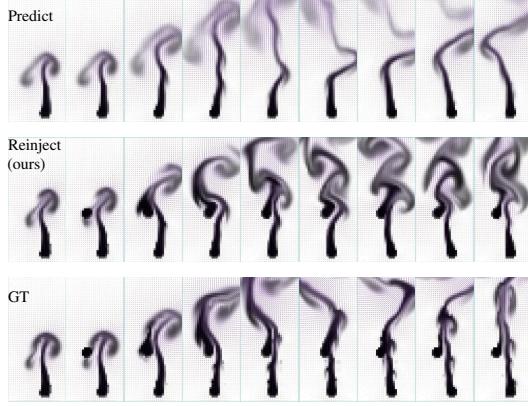wards learning a latent space mapping that is suitable for

Figure 6: An additional inflow, unseen in training, is placed during prediction. In contrast to our approach (center), the prediction only approach (top) is not able to capture the second inflow.

Table 2: Evaluations for *Vel* and *VelDen* predictions; rotating and moving cup: a) Internal predictions $n_i$; b) Simple split vs. LSS; c) LS dimensionality comparison

| | *Vel* | | *VelDen* | |
|---|---|---|---|---|
| a) $n_i$ | PSNR $u$ | PSNR $\rho$ | PSNR $u$ | PSNR $\rho$ |
| 1 | 33.04 | 25.19 | 33.11 | 22.28 |
| 6 | 35.89 | 26.28 | 36.07 | 25.41 |
| 12 | **36.61** | **26.61** | **36.61** | **25.69** |
| b) Type | PSNR $u$ | PSNR $\rho$ | PSNR $u$ | PSNR $\rho$ |
| Simple Split | 26.95 | 16.35 | **29.51** | 15.63 |
| LSS (ours) | **30.28** | **17.66** | 29.11 | **17.12** |
| c) LS Dimension $|c|$ | PSNR $u$ | PSNR $\rho$ | PSNR $u$ | PSNR $\rho$ |
| 16 | 30.28 | 17.66 | 29.11 | 17.12 |
| 32 | 30.40 | 17.64 | **31.58** | 18.90 |
| 48 | **30.86** | **17.92** | 30.97 | **18.96** |

temporal predictions. Intuitively, changes over time require relatively small and smooth updates, which results in the visually more regular curves shown in Figure 3c.

**Simple LS Division vs. LSS** A simple approach to arrive at a latent space with a clear subdivision in terms of input quantities is to use two separate spatial AEs for the individual input quantities. After encoding, the two resulting latent space points $c_{vel} = E_{vel}(u)$ and $c_{den} = E_{den}(\rho)$ can be concatenated, yielding $c_{simple} = [c_{vel}, c_{den}]$. In contrast to the simple approach, our LSS directly encodes both quantities with a single AE as $c_{LSS} = E([u, \rho])$ and enforces the subdivision with a soft-constraint. The combined training with the prediction network is performed identical for both spatial compression versions. It becomes apparent from the results in Table 2 b), that the network trained with our soft-constraint outperforms the simple split variant. Especially, when reinjecting the simulation density in the *Vel* benchmarks, we see a better PSNR value of 30.28 for $u$ and 17.66 for $\rho$ for our method in comparison to 26.95 and 16.35 for $u$ and $\rho$, respectively. The reason for this is that the simple split version can not take advantage of synergistic effects in the input data, since both input quantities are encoded in separate AEs. In contrast, our method uses the synergies of the physically interconnected velocity and density fields and robustly predicts future time steps.

**Internal Iterations** We compare the internal iteration count of the prediction network in the training process in Table 2 a). By performing multiple recurrent evaluations of our temporal network already in the training process we minimize the additive error build-up of many consecutive predictions. To fight the additive prediction errors over a long time horizon is important to arrive at a robust and exact predicted sequence. We chose the values of 1, 6 and 12 internal iterations for our comparison. It becomes apparent, that the network trained with 12 internal iterations and thereby the

longest prediction horizon is superior in both evaluations. It should be noted, that the predictions with reinjection of the physical density field (*Vel*) have a lower error on the density than the prediction-only (*VelDen*) approach, e.g. a PSNR value of 26.61 in contrast to 25.69 for the density field of the 12 iteration version. This supports the usefulness of our proposed latent space subdivision, that is needed to reinject external information.

**Latent Space Dimensionality** The latent space dimensionality has a major impact on the resulting weight count of the autoencoder as well as the complexity of the temporal predictions and thereby their difficulty. In the following we compare latent space sizes of 16, 32 and 48. When it comes to prediction only (*VelDen*), the PSNR is better for a larger latent space dimensionality. In contrast to this observation, the PSNR value is on the same level for all latent space sizes, when the simulation density is reinjected (*Vel*). For this reason we used a latent space dimensionality of 16 for all further comparisons. Due to bouyancy, the velocity and density of our smoke simulations are loosely coupled. Thus, additional weights do not increase the overall performance when the reconstruction of the individual parts of the respective input quantities already converged.

**Latent Space Subdivision vs. No-Split** To evaluate the usefulness of our LSS method that supports reinjection of external information, it is compared to a classic network setup that does not support that and instead performs a regular prediction. When only performing a regular prediction, the step-wise prediction errors accumulate. Without the reinjection of the externally driven density field into our prediction process, the quality of the outcome decreases drastically. Due to the bouyancy based coupling of density and velocity this effect intensifies.

In Table 3 we compare the long-term temporal prediction performance of a 0.0 (no-split) version and our 0.66 LSS version over a time horizon of 400 simulation steps. Those split numbers correspond to the percentage designated for

Table 3: LSS and no-split comparison; rotating and moving cup; 400 time steps

| LS Split | Type | PSNR $u$ | PSNR $\rho$ |
|---|---|---|---|
| 0.0 (no-split) | *VelDen* | 17.41 | 10.71 |
| 0.66 | *VelDen* | 26.81 | 17.07 |
| 0.66 | *Vel* | **28.15** | **21.74** |

the velocity part of the latent space, e.g., a value of 0.66 means that 66% of the latent space are used for the encoded velocity. Since the no-split version contains a classic AE without any latent space constraints, the evaluation can only be performed for the prediction-only (*VelDen*) benchmark, since the reinjection of density is not possible with the classic AE. Our LSS 0.66 version with a density PSNR value of 21.74 clearly outperforms the no-split version with a density PSNR value of 10.71. Due to the reinjection capabilities of our method, the resulting prediction remains stable whereas the classic no-split approach fails to capture the flow throughout the prediction horizon and even produces unphysical behavior (see Figure 4).

# 6. Results

We demonstrate the effectiveness of the subdivided latent space with several generalization tests. As shown in Figure 5, our method is capable of predicting the fluid motion even when an obstacle is placed in the domain. Due to our split latent space, the obstacle can be passively injected into the prediction process by supplying an encoded density field with a masked out obstacle region. We replaced the density part of the latent space with its encoded state after advecting it with our predicted velocity field. The prediction without
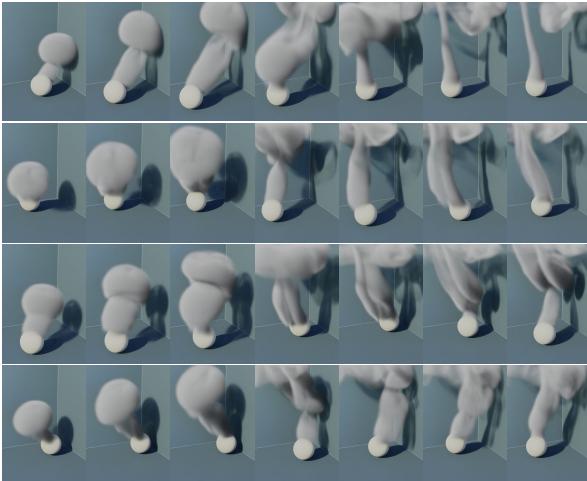


Figure 7: Four different 3D moving smoke predictions. The movement is completely predicted by our proposed method.

injection of external information is not capturing the obstacle and thereby deviates from the ground truth. In Figure 6 we show that our method is capable of predicting the fluid

motion even when a new inflow region is added that was not seen while training. Our network performs reasonably well in all tested experiments due to the reinjection capabilities provided by the latent space subdivision, even though the generalization scenes were not part of the training data.

Table 4: Average timing of a simulation step computed via regular pressure solve and our *Vel* prediction scheme. While the former scales with data complexity, ours scales linearly with the domain dimension. Average of 5 scenes with 100 simulation steps each. Measured on Intel(R) Xeon(R) E5-1650 v3 (3.50GHz) and Nvidia GeForce RTX 2070.

| Scene | Resolution | Type | Solve [s] | Total [s] |
|---|---|---|---|---|
| Rot. mov. cup 3D | $48^3$ | Simulation | 0.891 | 0.960 |
| Rot. mov. cup 3D | $48^3$ | Prediction | **0.074** | **0.156** |
| Mov. smoke 3D | $48^3$ | Simulation | 0.472 | 0.537 |
| Mov. smoke 3D | $48^3$ | Prediction | **0.059** | **0.132** |
| Rot. mov. cup | $64^2$ | Simulation | 0.041 | 0.044 |
| Rot. mov. cup | $64^2$ | Prediction | **0.012** | **0.019** |
| Rot. cup | $48^2$ | Simulation | 0.018 | 0.019 |
| Rot. cup | $48^2$ | Prediction | **0.011** | **0.015** |

To demonstrate the capabilities of our method, we trained a 3D version of our network on the moving smoke scene; selected frames are shown in Figure 7. Additionally, we compared the runtime performance of our networks to the regular solver that was used for generating the training data (see Table 4). Even though, we need to decode and encode the density field due to our reinjection method and thereby copy it from GPU to CPU memory, we still arrive at a performance measure of 0.059 seconds for an average prediction step in our 3D scene. For comparison, a traditional multi-threaded CPU-based solver takes 0.472 seconds on average for a simulation step for the same scenes.

# 7. Conclusion & Future Work

We have demonstrated an approach for subdividing latent spaces in a controlled manner, to improve generalization and long-term stability of physics predictions. In combination with our time-aware end-to-end training that learns a reduced representation together with the time prediction, this makes it possible to predict sequences with several hundred roll-out steps. In addition, our trained networks can be evaluated very efficiently, and yield significant speed-ups compared to traditional solvers. As future work, we believe it will be interesting to further extend the generalizing capabilities of our network such that it can cover a wider range of physical behavior. In addition, it will be interesting to explore different architectures to reduce the hardware requirements for training large 3D models with our approach.

## Acknowledgements

## References

Ando, R., Thürey, N., and Wojtan, C. A dimension-reduced pressure solver for liquid simulations. *Comput. Graph. Forum*, 34(2):473–480, May 2015. ISSN 0167-7055. doi: 10.1111/cgf.12576. URL http://dx.doi.org/10.1111/cgf.12576.

Bridson, R. *Fluid Simulation for Computer Graphics*. CRC Press, 2015.

Bromley, J., Guyon, I., LeCun, Y., Säckinger, E., and Shah, R. Signature verification using a "siamese" time delay neural network. In *Proceedings of the 6th International Conference on Neural Information Processing Systems*, NIPS'93, pp. 737–744, 1993.

Brunton, S., Noack, B., and Koumoutsakos, P. Machine Learning for Fluid Mechanics. may 2019. URL http://arxiv.org/abs/1905.11075.

Chu, M. and Thuerey, N. Data-driven synthesis of smoke flows with CNN-based feature descriptors. *ACM Trans. Graph.*, 36(4)(69), 2017.

Cui, Q., Sen, P., and Kim, T. Scalable laplacian eigenfluids. *ACM Trans. Graph.*, 37(4):87:1–87:12, July 2018. ISSN 0730-0301. doi: 10.1145/3197517.3201352. URL http://doi.acm.org/10.1145/3197517.3201352.

De Witt, T., Lessig, C., and Fiume, E. Fluid simulation using laplacian eigenfunctions. *ACM Trans. Graph.*, 31(1):10:1–10:11, February 2012. ISSN 0730-0301. doi: 10.1145/2077341.2077351. URL http://doi.acm.org/10.1145/2077341.2077351.

Gupta, M. and Narasimhan, S. G. Legendre fluids: A unified framework for analytic reduced space modeling and rendering of participating media. In *Proceedings of the 2007 ACM SIGGRAPH/Eurographics Symposium on Computer Animation*, SCA '07, pp. 17–25, Aire-la-Ville, Switzerland, Switzerland, 2007. Eurographics Association. ISBN 978-1-59593-624-0. URL http://dl.acm.org/citation.cfm?id=1272690.1272693.

Jones, A. D., Sen, P., and Kim, T. Compressing fluid subspaces. In *Symposium on Computer Animation*, pp. 77–84. ACM/Eurographics, 2016.

Kim, B., C. Azevedo, V., Thuerey, N., Kim, T., Gross, M., and Solenthaler, B. Deep Fluids: A Generative Network for Parameterized Fluid Simulations. *Computer Graphics Forum (Proc. Eurographics)*, 38(2), 2019.

Kim, T. and Delaney, J. Subspace fluid re-simulation. *ACM Trans. Graph.*, 32(4):62:1–62:9, July 2013. ISSN 0730-0301. doi: 10.1145/2461912.2461987. URL http://doi.acm.org/10.1145/2461912.2461987.

Ladický, L., Jeong, S., Solenthaler, B., Pollefeys, M., and Gross, M. Data-driven fluid simulations using regression forests. *ACM Transactions on Graphics*, 34(6):1–9, oct 2015.

Liu, B., Mason, G., Hodgson, J., Tong, Y., and Desbrun, M. Model-reduced variational fluid simulation. *ACM Trans. Graph.*, 34(6):244:1–244:12, October 2015. ISSN 0730-0301. doi: 10.1145/2816795.2818130. URL http://doi.acm.org/10.1145/2816795.2818130.

Long, B. and Reinhard, E. Real-time fluid simulation using discrete sine/cosine transforms. In *Proceedings of the 2009 Symposium on Interactive 3D Graphics and Games*, I3D '09, pp. 99–106, New York, NY, USA, 2009. ACM. ISBN 978-1-60558-429-4. doi: 10.1145/1507149.1507165. URL http://doi.acm.org/10.1145/1507149.1507165.

Manteaux, P.-L., Vimont, U., Wojtan, C., Rohmer, D., and Cani, M.-P. Space-time sculpting of liquid animation. In *Proceedings of the 9th International Conference on Motion in Games*, MIG '16, pp. 61–71, New York, NY, USA, 2016. ACM. ISBN 978-1-4503-4592-7. doi: 10.1145/2994258.2994261. URL http://doi.acm.org/10.1145/2994258.2994261.

Morton, J., Jameson, A., Kochenderfer, M. J., and Witherden, F. Deep dynamical modeling and control of unsteady fluid flows. In *Proceedings of Neural Information Processing Systems*, 2018.

Radford, A., Metz, L., and Chintala, S. Unsupervised representation learning with deep convolutional generative adversarial networks. *Proc. ICLR*, 2016.

Raveendran, K., Wojtan, C., Thuerey, N., and Turk, G. Blending liquids. *ACM Trans. Graph.*, 33(4):137:1–137:10, July 2014. ISSN 0730-0301. doi: 10.1145/2601097.2601126. URL http://doi.acm.org/10.1145/2601097.2601126.

Sato, S., Dobashi, Y., and Nishita, T. Editing fluid animation using flow interpolation. *ACM Trans. Graph.*, 37(5):173:1–173:12, September 2018. ISSN 0730-0301. doi: 10.1145/3213771. URL http://doi.acm.org/10.1145/3213771.

Selle, A., Fedkiw, R., Kim, B., Liu, Y., and Rossignac, J. An Unconditionally Stable MacCormack Method. *J. Sci. Comput.*, 35(2-3):350–371, June 2008.

Thuerey, N. Interpolations of smoke and liquid simulations. *ACM Trans. Graph.*, 36(1):3:1–3:16, September 2016. ISSN 0730-0301. doi: 10.1145/2956233. URL http://doi.acm.org/10.1145/2956233.

Thuerey, N. and Pfaff, T. MantaFlow, 2018.

Tompson, J., Schlachter, K., Sprechmann, P., and Perlin, K. Accelerating eulerian fluid simulation with convolutional networks. In *Proceedings of the 34th ICML Vol. 70*, pp. 3424–3433. JMLR. org, 2017.

Treuille, A., Lewis, A., and Popović, Z. Model reduction for real-time fluids. In *ACM SIGGRAPH 2006 Papers*, SIGGRAPH '06, pp. 826–834, New York, NY, USA, 2006. ACM. ISBN 1-59593-364-6. doi: 10.1145/1179352.1141962. URL http://doi.acm.org/10.1145/1179352.1141962.

Um, K., Hu, X., and Thuerey, N. Liquid splash modeling with neural networks. In *Computer Graphics Forum*, volume 37, pp. 171–182. Wiley Online Library, 2018.

Wicke, M., Stanton, M., and Treuille, A. Modular bases for fluid dynamics. *ACM Trans. Graph.*, 28(3):39, August 2009.

Wiewel, S., Becher, M., and Thuerey, N. Latent space physics: Towards learning the temporal evolution of fluid flow. *Computer Graphics Forum*, 38(2):71–82, 2019.

Wu, J., Zhang, C., Xue, T., Freeman, B., and Tenenbaum, J. Learning a probabilistic latent space of object shapes via 3d generative-adversarial modeling. In *Advances in Neural Information Processing Systems*, pp. 82–90, 2016.

Xie, Y., Franz, E., Chu, M., and Thuerey, N. tempoGAN: A Temporally Coherent, Volumetric GAN for Super-resolution Fluid Flow. *SIGGRAPH*, 2018.

Yang, C., Yang, X., and Xiao, X. Data-driven projection method in fluid simulation. *Computer Animation and Virtual Worlds*, 27(3-4):415–424, may 2016. ISSN 15464261. doi: 10.1002/cav.1695. URL http://doi.wiley.com/10.1002/cav.1695.

# Supplemental Document for Latent Space Subdivision: Stable and Controllable Time Predictions for Fluid Flow

## A. Evaluation

### A.1. Prediction Window Size

The prediction window $w$ describes the count of consecutive time steps that are taken as input by the temporal prediction network. In our comparison we tested window sizes ranging from 2 over 3 up to 4 consecutive input steps. The results in terms of PSNR values are displayed in Table 5 and Table 6.

Table 5: Prediction window $w$ comparison *Vel*

| $w$ | PSNR $u$ | PSNR $\rho$ |
|---|---|---|
| 4 | 29.67 | 17.04 |
| 3 | 29.79 | 16.87 |
| 2 | **30.28** | **17.66** |

Table 6: Prediction window $w$ comparison *VelDen*

| $w$ | PSNR $u$ | PSNR $\rho$ |
|---|---|---|
| 4 | **29.98** | **18.24** |
| 3 | 29.52 | 17.84 |
| 2 | 29.11 | 17.12 |

It becomes apparent that the prediction-only approach (*VelDen*) benefits from a larger input window, whereas the *Vel* approach with reinjected external information performs best with a smaller input window.

### A.2. Latent Space Split Percentage

We evaluated the impact of the latent space split percentage on three of our datasets. Therefore, we trained multiple models with different split percentages on the individual datasets. The comparison for our moving smoke scene is shown in Table 7 and Table 8. The latter are the results of the prediction-only evaluation (denoted *VelDen*), whereas the first table presents the results of our reinjected density approach (denoted *Vel*). In this experiment all split versions are outperformed by the no-split version in the prediction-only only setup with PSNR values of 29.71 and 18.03 for velocity and density, respectively.

Table 7: LS split comparison *Vel*; moving smoke

| LS Split | PSNR $u$ | PSNR $\rho$ |
|---|---|---|
| 0.33 | 28.07 | 15.84 |
| 0.5 | 29.28 | 16.49 |
| 0.66 | **30.28** | **17.66** |

Table 8: LS split comparison *VelDen*; moving smoke

| LS Split | PSNR $u$ | PSNR $\rho$ |
|---|---|---|
| 0.0 (no-split) | **29.71** | **18.03** |
| 0.33 | 28.49 | 16.63 |
| 0.5 | 29.06 | 17.38 |
| 0.66 | 29.11 | 17.12 |

Table 9: LS split comparison *Vel*; rotating cup

| LS Split | PSNR $u$ | PSNR $\rho$ |
|---|---|---|
| 0.33 | 36.67 | 28.46 |
| 0.5 | 36.66 | 29.22 |
| 0.66 | **38.52** | **29.73** |

Table 10: LS split comparison *VelDen*; rotating cup

| LS Split | PSNR $u$ | PSNR $\rho$ |
|---|---|---|
| 0.0 (no-split) | **37.90** | 22.68 |
| 0.33 | 37.52 | 25.01 |
| 0.5 | 36.77 | 25.13 |
| 0.66 | 37.57 | **25.32** |

Table 11: LS split comparison *Vel*; rotating and moving cup

| LS Split | PSNR $u$ | PSNR $\rho$ |
|---|---|---|
| 0.33 | 35.67 | 25.10 |
| 0.5 | **36.94** | **26.59** |
| 0.66 | 36.50 | 26.25 |

Table 12: LS split comparison *VelDen*; rotating and moving cup

| LS Split | PSNR $u$ | PSNR $\rho$ |
|---|---|---|
| 0.33 | **37.89** | **26.45** |
| 0.5 | 37.30 | 26.16 |
| 0.66 | 37.56 | 26.14 |

Table 13: LSS and no-split comparison; rotating cup; 100 time steps

| LS Split | Type | PSNR $u$ | PSNR $\rho$ |
|---|---|---|---|
| 0.0 (no-split) | *VelDen* | 37.90 | 22.68 |
| 0.66 | *VelDen* | 37.57 | 25.32 |
| 0.66 | *Vel* | **38.52** | **29.73** |

In contrast, the networks trained on the rotating cup dataset behave different as shown in Table 9 and Table 10. The classic no-split version is outperformed by all other split versions in terms of density PSNR values in the prediction-only (*VelDen*) setup. In the reinjected density evaluation (*Vel*), the benefit of latent space splitting becomes even more apparent when comparing the PSNR values of velocity, 38.52 and density, 29.73 of the 0.66 network with the velocity PSNR of 37.90 and density PSNR 22.68 of the no-split version.

### A.3. Latent Space Subdivision vs. No-Split

In this section we present additional results for our rotating cup dataset. See the main document for a long-term comparison of LSS vs. no-split for our more complicated moving and rotating cup dataset. In Table 13 we compare the temporal prediction performance of a 0.0 (no-split) version and our 0.66 LSS version over a time horizon of 100 simulation steps. Our LSS 0.66 version with a density PSNR value of 29.73 clearly outperforms the no-split version with a density PSNR value of 22.68.

### A.4. Generalization

Additionally, we show in Figure 8 that our method recovers from the removal of smoke in a certain sink-region and is capable of predicting the fluid motion.
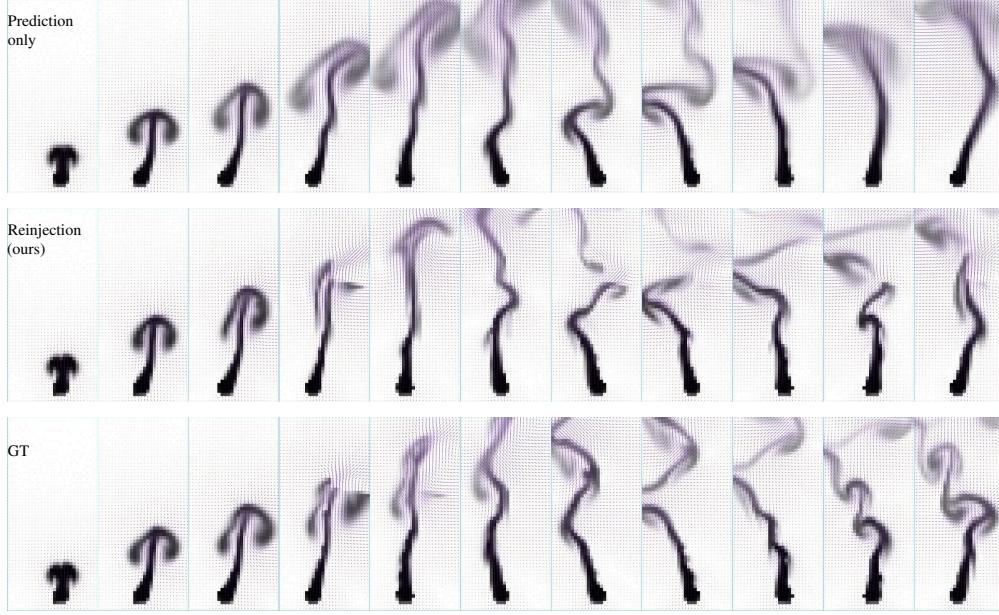
Figure 8: An sink is placed in the upper right of our moving smoke scene. This was unseen during training. The prediction by our proposed method remains stable and realistic. In the second row density reinjection was applied. In the top row no external information was injected. Thus, the sink can't be processed by the network.

## B. Fluid Flow Data

Our work concentrates on single-phase flows, modelled by a pressure-velocity formulation of the incompressible Navier-Stokes equations as highlighted in Section 2. Thereby, we apply a classic NS solver to simulate our smoke flows based on R. Bridson (2015). In addition to Section 4, more information about the simulation procedure is provided in the following.

### B.1. Simulation Setup

The linear system for pressure projection is solved with a conjugate gradient method. The conjugate gradient (CG) solver accuracy is set to $1 \cdot 10^{-4}$ for our moving smoke dataset, whereas an accuracy of $1 \cdot 10^{-3}$ is utilized for the moving cup datasets. We generated all our datasets with a time step of $0.5$. Depending on the behavioral requirements of our different experiments with rising, hot and sinking, cold smoke we use the Boussinesq model with the smoke density in combination with a gravity constant of $(0.0, -4 \cdot 10^{-3}, 0.0)$ for the moving and rising smoke and $(0.0, 1 \cdot 10^{-3}, 0.0)$ for the rotating cup dataset. To arrive at a more turbulent flow behavior, the gravity constant was set to $(0.0, 1 \cdot 10^{-2}, 0.0)$ for our moving and rotating cup dataset. We do not apply other forces or additional viscosity. We purely rely on numerical diffusion to introduce viscosity effects.

In combination with the quantities required by our classic NS setup, namely flow velocity $u$, pressure $p$ and density $\rho$, we also need a flag grid $f$, an obstacle velocity field $u_{obs}$ and the corresponding obstacle levelset for our obstacle supporting scenes. Thereby our density $\rho$ is passively advected within the flow velocity $u$.

To handle the obstacle movement accordingly, we calculate the obstacle velocity field by evaluating the displacement per mesh vertex of the previous to the current time step and applying the interpolated velocities to the according grid cells of the obstacle velocity field. Afterwards, the obstacle velocity field values are averaged to represent a correct discretization.

In Algorithm 1 the simulation procedure of the moving smoke dataset is shown. For our obstacle datasets the procedure in Algorithm 2 is used, with the prediction algorithm given in Algorithm 3. Boundary conditions are abbreviated with BC in these algorithm.

---

**Algorithm 1** Moving smoke simulation

> **while** $t \to t + 1$ **do**
> $\quad \rho \leftarrow \text{applyInflowSource}(\rho, s)$
> $\quad \rho \leftarrow \text{advect}(\rho, u)$
> $\quad u \leftarrow \text{advect}(u, u)$
> $\quad f \leftarrow \text{setWallBCs}(f, u)$
> $\quad u \leftarrow \text{addBuoyancy}(\rho, u, f, g)$
> $\quad p \leftarrow \text{solvePressure}(f, u)$
> $\quad u \leftarrow \text{correctVelocity}(u, p)$
> **end while**

---

**Algorithm 2** Rotating and moving cup

1:  **while** $t \rightarrow t+1$ **do**
2:     $\rho \leftarrow$ applyInflowSource($\rho$, $s$)
3:     $\rho \leftarrow$ advect($\rho$, $\boldsymbol{u}$)
4:     $\boldsymbol{u} \leftarrow$ advect($\boldsymbol{u}$, $\boldsymbol{u}$)
5:     $\boldsymbol{u}_{obs} \leftarrow$ computeObstacleVelocity($obstacle^t$, $obstacle^{t+1}$)
6:     $f \leftarrow$ setObstacleFlags($obstacle^t$)
7:     $f \leftarrow$ setWallBCs($f$, $\boldsymbol{u}$, $obstacle^t$, $\boldsymbol{u}_{obs}$)
8:     $\boldsymbol{u} \leftarrow$ addBuoyancy($\rho$, $\boldsymbol{u}$, $f$, $\boldsymbol{g}$)
9:     $p \leftarrow$ solvePressure($f$, $\boldsymbol{u}$)
10:    $\boldsymbol{u} \leftarrow$ correctVelocity($\boldsymbol{u}$, $p$)
11: **end while**

**Algorithm 3** Rotating and moving cup network prediction *Vel*

1:  **while** $t \rightarrow t+1$ **do**
2:     $\rho \leftarrow$ applyInflowSource($\rho$, $s$)
3:     $\rho \leftarrow$ advect($\rho$, $\boldsymbol{u}$)
4:     $\boldsymbol{u} \leftarrow$ advect($\boldsymbol{u}$, $\boldsymbol{u}$)
5:     $\boldsymbol{u}_{obs} \leftarrow$ computeObstacleVelocity($obstacle^t$, $obstacle^{t+1}$)
6:     $f \leftarrow$ setObstacleFlags($obstacle^t$)
7:     $f \leftarrow$ setWallBCs($f$, $\boldsymbol{u}$, $obstacle^t$, $\boldsymbol{u}_{obs}$)
8:     $\dot{\boldsymbol{c}}^t \leftarrow$ encode($\tilde{\boldsymbol{u}}^t$, $\rho^t$)
9:     $\hat{\boldsymbol{c}}^t \leftarrow [\tilde{\boldsymbol{c}}^t_{vel}, \dot{\boldsymbol{c}}^t_{den}]$
10:    $\tilde{\boldsymbol{c}}^{t+1} \leftarrow$ predict($\hat{\boldsymbol{c}}^{t-1}$, $\hat{\boldsymbol{c}}^t$)
11:    $\tilde{\boldsymbol{u}}^{t+1}, \tilde{\rho}^{t+1} \leftarrow$ decode($\tilde{\boldsymbol{c}}^{t+1}$){} $\tilde{\rho}^{t+1}$ is not used
12:    $\boldsymbol{u}^{t+1} \leftarrow \tilde{\boldsymbol{u}}^{t+1}$ {o}verwrite the velocity with the prediction
13: **end while**

### B.2. Training Datasets

In the following multiple simulations contained in our training data set are displayed.



Figure 9: Four example sequences of our moving smoke dataset. For visualization purposes we display frames 20 to 200 with a step size of 20 for the respective scenes. The smoke density is shown as black.
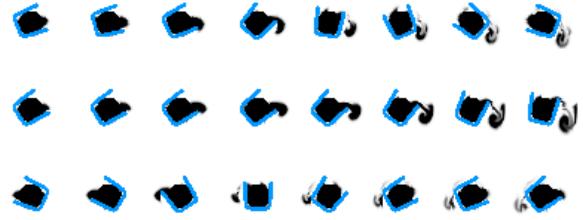


Figure 10: Three example sequences of our rotating cup dataset. For visualization purposes we display frames 40 to 180 with a step size of 20 for the respective scenes. The cup-shaped obstacle is highlighted in blue, whereas the smoke density is shown as black.



Figure 11: Three example sequences of our rotating and moving cup dataset. For visualization purposes we display frames 40 to 180 with a step size of 20 for the respective scenes. The cup-shaped obstacle is highlighted in blue, whereas the smoke density is shown as black.