

Motion Priors for Pose Estimation and Animation Workflows

Luke Smith

Master Thesis
May 2023

Prof. Dr. Robert W. Sumner



Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich



computer graphics laboratory

Abstract

Motion capture, the act of automatically recording the motion of a person or object, has solidified it's position as a key technology in modern animation pipelines due to its numerous benefits. It provides real time, high quality, physically plausible results, however the upfront costs for tailored software and hardware, cumbersome motion capture suits, and other issues leads to a high barrier not only for individual artists but also small teams. At Disney Research|Studios, this barrier is being tackled with a system that aims to capture motion directly from RGB video. The system performs pose estimation on the individual frames, and runs the result through an optimizer to clean up the motion. While promising, the system occasionally struggles with certain artifacts and with some occluded motions. The goal of this thesis is therefore to consider how we might learn a general human motion model that could improve upon these shortcomings. To begin with we investigate a state of the art approach based on Variational Auto Encoders, and attempt to improve upon its issues to make it useable for our purposes. We then consider a relatively new class of models, that of diffusion models, which is proving to be a powerful tool in many domains. We consider the application of such models in the context of motion modelling, and implement and evaluate a baseline model with the view of coming to a go/no-go decision on the feasibility of such models for our purposes.

Zusammenfassung

Motion Capture, die automatische Aufzeichnung der Bewegung einer Person oder eines Objekts, hat sich aufgrund ihrer zahlreichen Vorteile als Schlüsseltechnologie in modernen Animationspipelines etabliert. Sie liefert qualitativ hochwertige, physikalisch plausible Ergebnisse in Echtzeit, doch die Vorabkosten für maßgeschneiderte Software und Hardware, umständliche Motion-Capture-Anzüge und andere Probleme führen zu einer hohen Hürde nicht nur für einzelne Künstler, sondern auch für kleine Teams. Bei Disney ResearchStudios wird diese Hürde mit einem System angegangen, das darauf abzielt, Bewegungen direkt aus RGB-Videos zu erfassen. Das System führt eine Posenschätzung für die einzelnen Frames durch und lässt das Ergebnis durch einen Optimierer laufen, um die Bewegung zu bereinigen. Das System ist zwar vielversprechend, hat aber gelegentlich Probleme mit bestimmten Artefakten und mit einigen verdeckten Bewegungen. Ziel dieser Arbeit ist es daher zu untersuchen, wie ein allgemeines menschliches Bewegungsmodell erlernt werden kann, das diese Unzulänglichkeiten beseitigt. Zunächst untersuchen wir einen aktuellen Ansatz, der auf Variational Auto Encoders basiert, und versuchen, seine Probleme zu verbessern, um ihn für unsere Zwecke nutzbar zu machen. Anschließend betrachten wir eine relativ neue Klasse von Modellen, die Diffusionsmodelle, die sich in vielen Bereichen als leistungsfähiges Werkzeug erweisen. Wir betrachten die Anwendung solcher Modelle im Zusammenhang mit der Bewegungsmodellierung und implementieren und bewerten ein Basismodell, um eine Entscheidung über die Durchführbarkeit solcher Modelle für unsere Zwecke treffen zu können.

Contents

List of Figures

Introduction

Motion capture is an integral part of many modern animation pipelines. It can be defined simply as the act of recording, by automatic means, the motion of a person, animal or object. Throughout this thesis we will primarily interest ourselves in the capturing of human motion. The advantages of motion capture, over traditional animation, are numerous and important in their scope. Motion capture allows for

- quasi-real time results
- high quality motion with realistic object interactions
- often reduced costs as compared to hand animation
- complexity that is constant with respect to the motion being captured

however also has it's inherent costs, notably

- upfront costs for tailored software and hardware
- cumbersome motion capture suits
- lengthy setup times and a non-trivial capture systems requiring expert knowledge
- artifacts due to retargeting of skeletons.

These drawbacks mean that the technology is often prohibitive for small teams and individual artists/animations. It is clear though that overcoming these barriers would provide huge benefits and open up the technology to a wide range of new users. For example individual artists could more efficiently prototype and develop animation sequences starting from a self captured motion, and small teams could readily make use of this technology to do more with less.

The important question therefore is simply; **how might we create a motion capture system without such upfront costs, specialized hardware and need for motion capture suits?** At Disney ResearchStudios an approach is being investigated that aims to capture motion directly

1 Introduction

from RGB video. The system currently follows the following steps:

1. record a video of a human motion sequence (potentially from multiple angles)
2. perform per frame pose estimation using a machine learning model
3. run the resultant motion sequence through an optimiser to improve the quality of the motion (and potentially lift to 3d)

and shows great promise, however has a number of drawbacks. The most notable drawback comes from the fact that the per frame pose estimation system does not produce temporally consistent results. The predictions are made independantly per frame, hence the result motion is jittery, and can contain other artifacts. This is the motivation for the introduction of an optimisation system at the end of the pipeline to counteract these artifacts. However it must be noted that, although smoothness is mostly acheived, it does not fix all issues. The optimiser struggles to handle occluded motion sequences, it does not always deal with limb flips (where the left/right leg/arm are predictions are horizontally flipped for a frame), and cannot always fix completely wrong but confident predictions. It is of interest therefore to try and improve this aspect of the pipeline.

The data driven approach proposed in this thesis, in its most abstract form, is simply that of learning a model that understands human motion. Depending on its conception, such a model could be employed in a number of manners; it could be directly applied to the task of rectifying a motion sequence, it could be used as a loss in the existing optimiser, or it could even simply be used to improve upon certain failings of the optimizer as an additional step to the pipeline, such as fixing occluded motion.

This thesis tackles the task of exploring such motion models with a primary goal of shedding light on the most promising model architecture and the most effictive manner in which such a model might be used. We consider a state of the art approach in Chapter ??, evaluate it in Section ?? and attempt to improve upon it in Section ?? . This does not prove as fruitful as hoped and so we move on to a new class of powerful models called diffusion models in Chapter ?? . We explore the use of these models in the context of motion modelling in Section ?? and finally evaluate them on a variety of tasks with a view of better understanding their capabilities in Section ?? .

Related Work

The study of synthesising human motion has a long history motivated in no small part by the desire to create realistic and captivating media in the gaming and film industries. Early adaptive methods rely on motion matching [?] [?] in which interpolation between similar motions from a database of captured motion is performed. This however does not scale well to out of database motion and often generates generic, non stylized motion, though efforts have been made to introduce learned aspects to such methods [?] to improve upon their shortcomings.

TODO: MORE?

More recent branches of motion modelling commonly base themselves upon the use of machine learning techniques, notably deep learning, to learn a prior over plausible motion. This is a more general approach that can be applied to a wider range of tasks, and shows promise in overcoming some of the issues of motion inbetweening, as such systems can learn to better generalise to out of distribution motion sequences.

Within the area of deep learning, many techniques have been investigated, temporal convolutions [?], recurrent models [?], and reinforcement learning [?] are but a few examples.

2.1 Motion AutoEncoders

A well explored model is that of the AutoEncoder (AE) [?] or Variation-AutoEncoder (VAE) [?]. These are popular models as they encourage the learning of a latent representation [?] of human motion, thus the intuition is that they learn not just to reproduce the data, but actually how humans move, providing a more robust prior.

TODO: MORE? c.f MEVA for nice related work section

Holden et al. [?] [?] present simple CNN based autoencoder architectures that operates on

2 Related Work

motion sequences. The notion of skeletal aware convolutions and pooling/unpooling operations for a VAE, alongside a sliding window method for motion rectification are presented by the authors of [?]. [?] presents a novel approach of leaning a latent space, then projecting directly to this latent space from a motion sequence using a separate model, again operating directly on a motion sequence. The authors of MEVA [?] postulate that a VAE often learns only smooth motion, as we are asking too much of the model, thus present a pipeline in which a smooth motion and coarse motion VAEs are jointly used. Holden et. al make another appearance with DeepPhase [?], an autoencoder with a latent space enforced to match sinusoidal functions that represent periodic motion. Contrary to the common trope of sequence level models, the authors of [?] and of [?] operate in a frame to frame regime, predicting temporally local change of motion. Finally, a number of works present the Conditional-VAE architecture [?] as a base with varying state representations, conditioning variables and loss terms, [?, ?, ?, ?].

As we can see the literature is rich and diverse, but we found ourselves drawn to the HuMoR model [?], due to it's state of the art performance and it's described use for the exact task that we desire to solve, that of rectifying a motion sequence captured through frame by frame pose estimation. The authors of [?] present a C-VAE architecture that learns a distribution over latent transitions, conditioned on the previous pose. They use this architecture alongside an optimisation method that rectifies human motion obtained from, among other modalities, RGB video through frame-by-frame pose estimation.

2.2 Motion Diffusion

With the notable success of diffusion models in the image generation literature [?, ?, ?], diffusion models have begun to spread into many other fields within machine learning [?], including to the the field of human motion modelling.

The authors of Avatars grow legs [?] denoise a sequence of SMPL [?] parameters condition on sparse tracking inputs, notably taking the form of the orientation/translation of a headset and two hand controllers. They show that plausible motion can be generate from very sparse signals, thus indicating to us that the use of diffusion models in the rectification of occluded motion sequences is promising. Next, PhysDiff [?] provides a text conditioned diffusion model with the unique use of a physics based motion projection step in the difussion process that helps to ensure physical plausibility of the generated motion. The authors of EDGE [?] propose an attention based, audio conditioned diffusion framework for dance motion generation. With a similar architecture to that of EDGE [?] but using transformers as the base of the denoising network such that the attention mechanism can more easily be exploited across the whole motion sequence, the authors of MDM [?] describe a text conditioned denoising architecture.

As we have seen, diffusion models can be employed in the field of motion modelling for a wide variety of tasks. Motion can be generated [?, ?, ?] conditioned on various inputs. Motion sequences can also be edited through inpainting [?, ?] to change only parts of a sequence, and motion inbetweening can also be acheived in a similar manner [?]. This wide variety of tasks that can be completed by a single model is a very attractive property of the diffusion framework, and inspired us to investigate these models later in the thesis.

HuMoR

The authors of HuMoR [?] present a novel approach for learning and using a plausible motion prior. They train a conditional VAE that learns a distribution over latent transitions, in a canonical reference frame, between *states* that consist of a root translation, 3D joint positions, joint angles, and the respective velocities. They most notably use this model as a prior in a 'test time optimisation', which generates plausible sequence motions optimising for an initial state and a sequence of transitions starting from frame by frame estimates (2D/3D joints or points clouds). This optimisation includes, alongside others, a motion prior term based upon the conditional distribution $p(z_t|x_{t-1})$ that encourages plausible motion for the learned sequence. Note that the CVAE decoder also predicts foot contacts alongside change in state, which are used as regularisers during their main use case of 'test time optimisation' (TestOps). The TestOps can operate on many modalities, 2D/3D joints, point clouds, etc., as the optimisation loss contains a Data Term ϵ_{data} that can be tailored to the modality as the HuMoR state is information rich, containing 3D joints (hence can fit to 2D joints through projection or directly to 3D) and can parametrise the SMPL model (hence the SMPL mesh can be correlated to point clouds).

The performance of HuMoR as described in the paper [?], alongside it's use in a problem that directly matches our own, lead us to evaluate and investigate the model Section ??, and subsequently try to extend and improve ?? upon it once the limitations became clear.

3.1 Relevant Model Details

The HuMoR architecture is that of a C-VAE [?], as can be seen in Figure ?? . The goal of this architecture is to learn a mapping from a given pose (

3.1.1 Method

Our investigation began with a largely qualitative evaluation of the HuMoR model which had two main aims. First was to stress test the system, to see where it failed, where it succeeded, and if the mentioned benefits in the HuMoR paper [?] were as described. Second was to evaluate the model with the defects of the current Disney Research|Studios system in mind to see if it could complement it's functionality, notably if it could improve upon occluded motion, joints flipping and confident but false predictions.

To achieve this goal, a selection of videos were taken in the Disney Research|Studios lab containing a variety of motion; fast, slow, and abnormal, and a including number of occluded scenes. The HuMoR system was ran on these videos and the results were investigated. The TestOps is performed in 3 stages, where only the 3rd makes use of the HuMoR motion model, we could therefore compare the stage 2 results to the stage 3 results to see where the model was providing an improvement over a more classical optimisation that includes smoothness and pose prior losses.

3.1.2 Advantages of HuMoR

We see a number of situations in which the model shows a clear improvement over the stage where the HuMoR model is not used.

In an occluded situation, as shown in Figure ??, where the 2d pose predictions don't have any information on the legs, the model manages to produce a realistic sitting motion.



Figure 3.1: HuMoR achieving occluded sitting

We also note in many of the less difficult videos, that HuMoR produces clean motion and deals with many abnormal movements without obvious issues. It therefore doesn't seem to regress the easy situations, but can improve the more difficult situations, notably occlusions.

3.1.3 Drawbacks of HuMoR

We noted that while HuMoR manages to sit when occluded, it often fails to walk. This seems to be due to the fact that OpenPose often predicts both legs on the frames just before the occlusions where only one leg is actually un-occluded, as can be seen in Figure ???. This results in a sequence of poses where the frames before an occlusion indicate that the person is no longer walking, hence making it significantly more difficult for HuMoR to begin walking again during

3.1 Relevant Model Details

the occlusion. This issue is thus probably largely due to HuMoRs' TestOps' dependence on OpenPose and thus it's inheritance of OpenPoses' failure points.

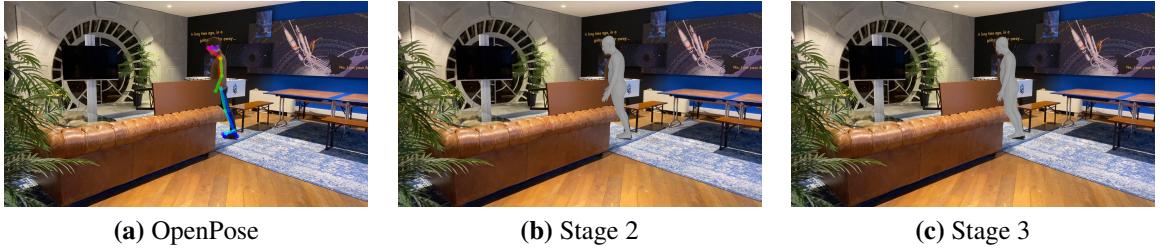


Figure 3.2: Occluded walking failure point, 2 legs predicted instead of 1

We also found that the choice of axis angle representation for various rotations led to the emergence of common known issue that arises from the discontinuities present in the representation [?]. The shortest path between certain angles in axis angle space can lead to a 360 degree rotation, as seen in Figure ??, among other issues.



Figure 3.3: Axis angle issue: dubious rotations

We found that the TestOps system would occasionally get itself into a tangle, as seen in Figure ??, we assume due to the autoregressive nature resulting in a potentially catastrophic cumulation of errors, and found this happened most notably for a sequence containing someone rolling on the floor. It is interesting to note that these sorts of motions were not often, if ever, present in the training data.



Figure 3.4: HuMoR in a tangle

Next, we saw overly smoothed motion for certain clips, most notably for particularly stylised motion such as dancing. We also found that if there is a single frame without Open Pose predictions then as of that frame the TestOps fails. Finally, and most importantly, we found that

the TestOps was extremely slow. It took around 20mins per 2s clip, and that we could batch at most 4 2s clips together, so 20mins per 8s on a Nvidia GeForce 1080ti with 11Gbs memory.

We concluded that there were a number of potential benefits, notably in occluded situations, and that many of the issues, including the axis angle representation issue, the dependence on humor, and smooth motion, have a good chance of being fixed with different design choices. The glaring issue however was the unreasonable amount of time the system took, rendering it unusable for our purposes as we wish the system to be close to real time.

3.1.4 Profiling

To investigate this speed issue, we profiled the code, as can be seen in Figure ???. We found that the program spent 90% of its time in the Stage 3 optimiser closure, 56% of its time performing the backwards step and 32% of its time in the rollout function. Hence it was clear that the speed issue was due to the slow act of rolling out and the large computation necessary to perform the backwards step on the large computation graph resulting from the rollout.

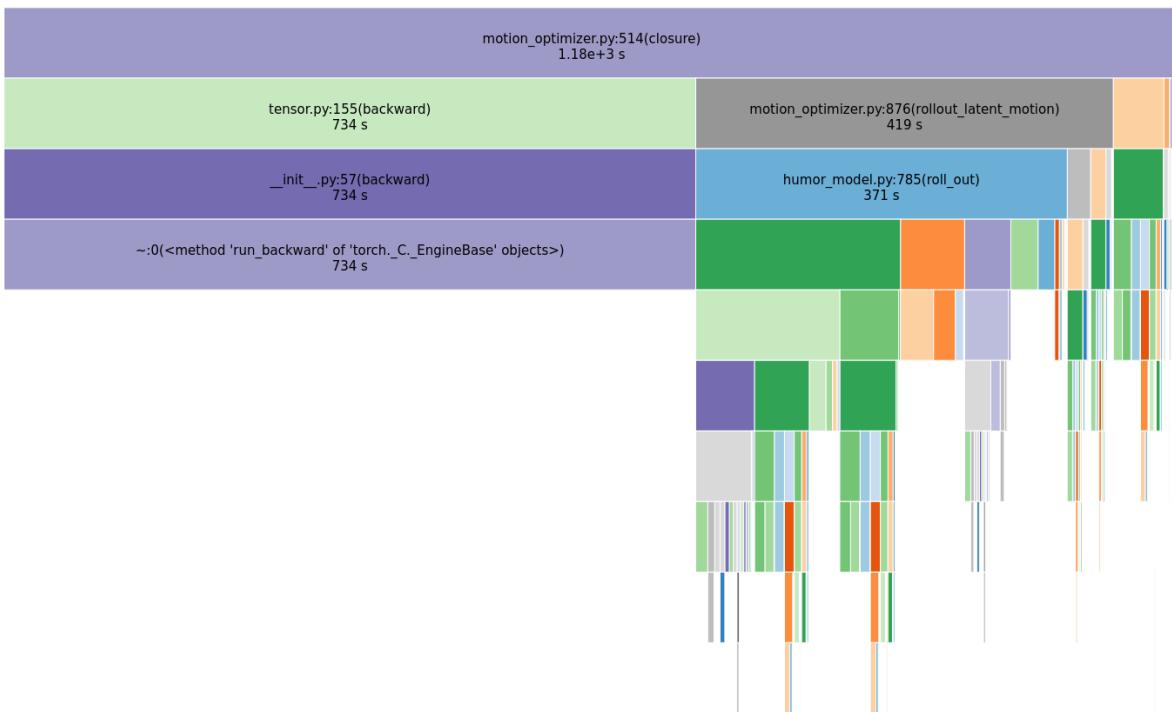


Figure 3.5: TestOps profiling

3.1.5 Conclusion

We came to the final conclusion that the HuMoR motion model alongside the TestOps optimisation showed potential, but that it must be sped up before we can conclusively say if it is of particular use to us. We note that to speedup the method, we must primarily focus on im-

proving/removing the rollout. Thus, in the next section, an attempt to speed up the TestOps is made.

3.2 Improving HuMoR TestOps

This section presents the attempt that was made at modify the TestOps of the HuMoR paper [?] in such a way that we attain comparable results but in a significantly shorter timespan.

TODO: See if I have introduced the terms properly (decoded sequence, x 's, x' 's, etc.)

3.2.1 Speeding up the TestOps

Realising that the main speed limitations are due to the concept of rollout over the entire sequence, we decided to break this long range dependence. Through short overlapping rollouts, it is possible to achieve more parallelisation and smaller computation graphs, and our intuition suggests that the need for such a large context window is not necessary, that only a small number of autoregressed frames are needed to be able to propagate information sensibly and to get meaningful gradients. Anecdotally, it would seem intuitive that if a video contains sitting, then walking, then sitting again, the second act of sitting would not depend on the early act of sitting, hence rolling out over the whole sequence seems unnecessary.

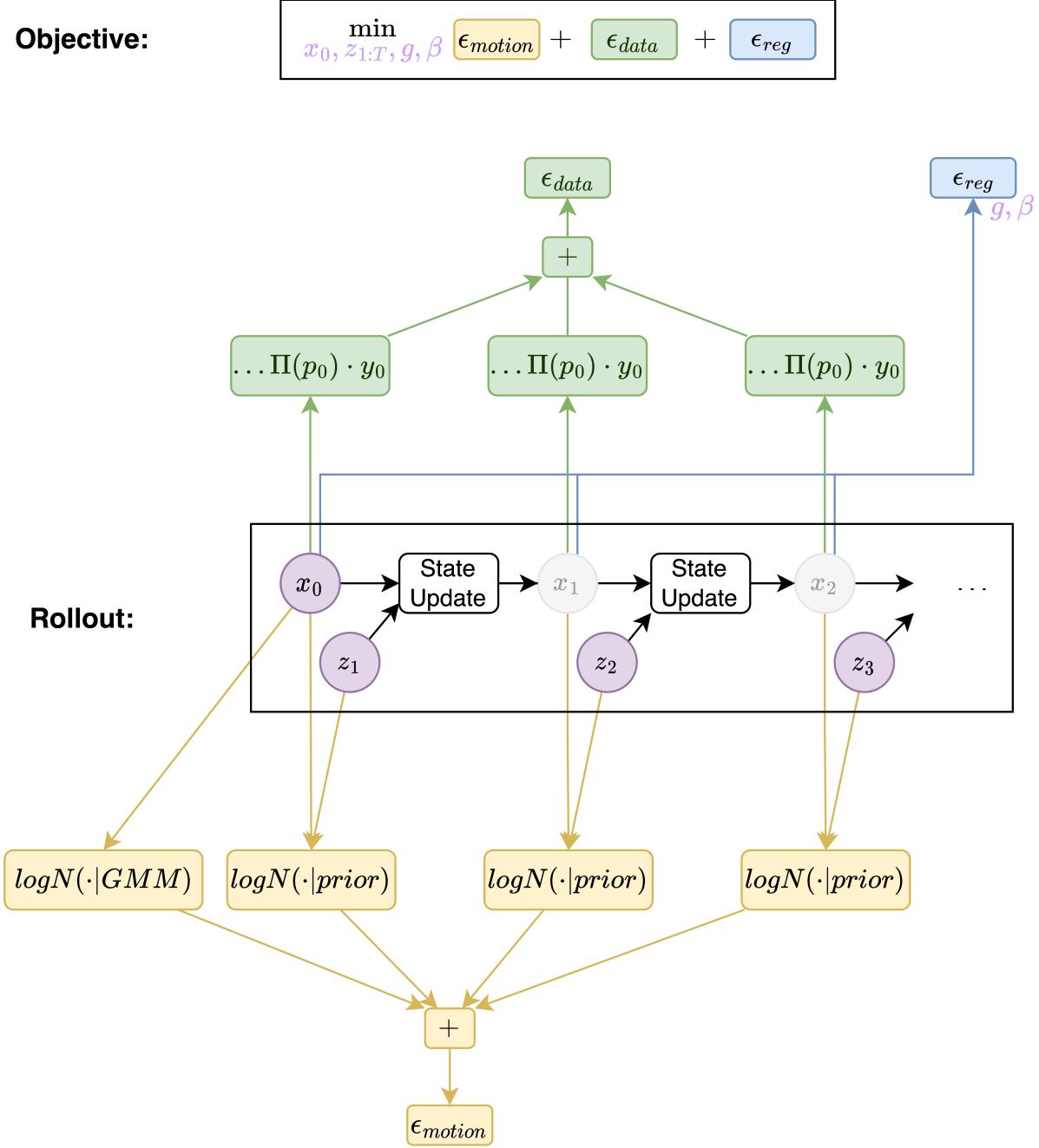
The HuMoR TestOps rollout is presented graphically in Figure ???. x is the state, and z are the latent variables, the optimised variables are highlighted in pink, and the losses are coloured, gray x 's indicate they are calculated during the optimisation from the optimised variables. As can be seen, first all the states are autoregressed from the initial state and the latent variables, and then the losses are applied. This can be seen to result in very long range dependencies in the computation graph, thus making the gradients time consuming to calculate.

TODO: Update Figure ?? and remove the 'NEW' norm, labels on green loss, and potentially red loss?

The most obvious way to break this autoregression is to maintain and optimise a separate sequence of x 's, updating this separate sequence with reference to the x' 's decoded from the latent variables as in Figure ???. We can take into account the decoded x' 's in several manners:

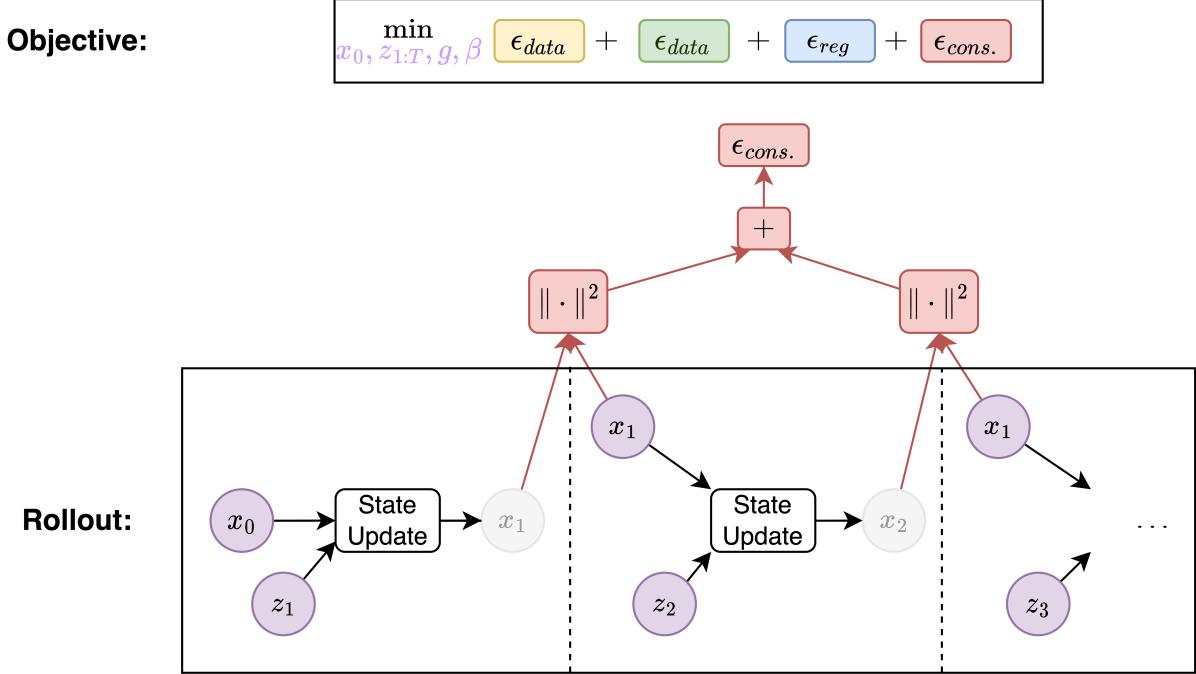
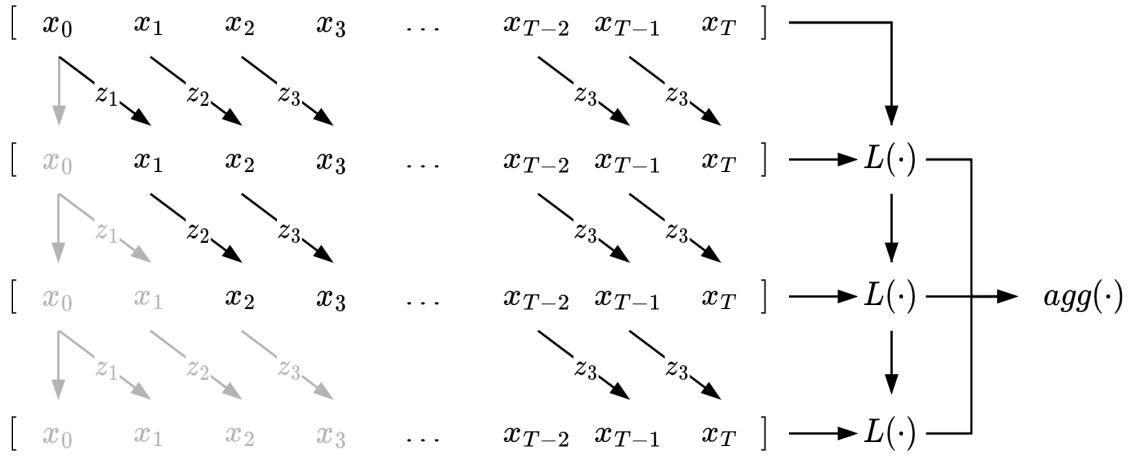
- Copy over
 - Directly replace the optimised sequence of x 's with the decoded x' 's.
- Blend
 - Perform a weighted addition of the optimised x 's and decoded x' 's.
- Loss term
 - Add a loss term on the difference between the optimised x 's and decoded x' 's.

All decoding steps can now be performed in parallel, rather than sequentially, which greatly reduces the computation time.

**Figure 3.6:** TestOps Computation Graph

This new method also allows us to experiment with different amounts of rollout. We can once again decode the decoded sequence of x 's to get x'' , and so on and so forth. These extra sequences are all decoded with the same latent variables thus the updates to the latent variables will take into account the losses on all the decoded sequences. Graphically this can be seen in Figure ??, where each next line . **TODO: describe the figure better**

The information contained in each x can flow forwards through the short range rollouts, and over the iterations of the optimiser. Note however that the information flow will be slower than

**Figure 3.7: Decoupled Computation Graph****Figure 3.8: Decoded sequences**

the HuMoR TestOps, as in the TestOps all subsequent x 's are regressed from the initial x_0 . In the HuMoR TestOps information can flow backwards from any future frame's x in one update through the gradients, as all x 's are linked in the computation graph through the auto-regression (though our intuition suggests that the information does not actually need to flow so far, as previously discussed). In the new method however, the information flow will be much more limited, as the rollouts will be significantly shorter, though again the information should flow through the iterations of the optimiser.

3.2.2 Implementation notes

A number of issues were encountered during the implementation that are worth mentioning. After implementing the new method, attempts were made to get the optimiser to produce sensible results. The general approach was to use a small number of losses, so as to get a better intuition of each loss, with the goal of better understanding how to balance the losses in the optimiser, but some errors in individual losses were encountered.

Firstly, as expected, the Axis-Angle representation became a problem. When running through functions to convert to root-local reference frame for the HuMoR model for decoding (as it operates on said reference frame), and then back to world reference frame for the opftimiser, the Axis-Angle vector flipped direction. Though the flipped vector still represented the same rotation, the angle representations of the same angle in the optimised sequence of x 's and the decoded sequence of x' 's were different, hence encorporating the information from decoded sequence cause the angle to move somewhere between the two flipped values which was no longer a sensible angle representation, thus resulting in root flips and other strange effects.

Secondly, an issue due to the SMPL model was found. When using a 2d reprojection loss (comparing the projected SMPL joints to the OpenPose 2d predictions) on a sequence where only the left arm and face had any OpenPose predictions, the legs were being moved by the optimiser. This was wrong as there were only losses on the arm and face joints, hence no gradients should flow to the legs. It was eventually found that the skinning weights of the SMPL model [?] contain spurious long range connections. 3% of the LBS (linear blend skinning) weights in the SMPL model are non-zero but less than $1e - 2$, which seems rather too low to have any meaningful effect on the skinning, but which allows for spurious gradients to flow.

For the comparison to the OpenPose predictions the OpenPose skeleton must be obtained from the SMPL mesh. Certain joints are regressed, other are simply taken directly as vertices from the mesh, as in the case of the nose joint in the OpenPose skeleton which is taken to be vertex 332 [?]. We noted that this vertex 332 is skinned 99.8% by the 'head' joint, but also 0.2% by the right hip, left knee and right knee. This issue was solved by pruning all weights below $1e-2$. It is interesting to note that there are known spurious connections in the SMPL blend shapes [?], but that we have found no reference to spurious connections in the LBS weights.

3.2.3 Experiments

TODO: explain why the initial rollout was bad, and the compounding of errors forward
TODO: maybe we should discuss more in the humor investagtion section about the bad initial rollout
TODO: Check everything is explicitly named, e.g have a graphic clearly naming the optimised x 's and the decoded x 's, and what the z 's are

To begin with, we looked into the initial z 's that were obtained from projecting the Stage 2 results through the HuMoR encoder. We found that in occluded situations these z 's encode a sitting motion without any optimisation as seen in Figure ??, that simply the projection leads them to encode the necessary occluded motion. However we note that when autoregressively rolled out over a longer sequence the results, in Figure ??, deviate largely from the initial sequence of x 's seen in Figure ?? . This shows us that though the initial z 's encode some sensible motion locally, and can be used to accurately recover the next frame, when they are

chained together small deviations from the x motion accumulate to create an overall deviation. For example if the arm is moved slightly too much by a single z , then all the future frames will have the arm in slightly the wrong position, and multiple z 's are wrong in the same direction, then the arm will raise up. This deviating sequence is the initial starting rollout in the Stage 3 of the HuMoR TestOps in which the rollout is optimised, though it seems unfortunate to us that from a sensible sequence of x 's obtained from Stage 2 we get such a bad sequence through the autoregressive rollout. We do however note that though while this starting point represents a sequence that deviates largely, the z 's might not need to move so far to avoid the accumulation effect. The intuition in our case was that we are starting from some sensible x 's and z 's, and therefore that it should be possible to refine the x 's by take into account the HuMoR model through the z 's whilst updating the z 's.



Figure 3.9: Stage 2 z 's encode a sitting motion

With some preliminary experiments in which we began optimising x 's and z 's we found that the most successful approach to blending the decoded and optimised x 's was to have an L1 loss on the distance between them. We saw that when simply copying over, the propagation of small errors forward was too strong (as described in Figure ??) and the system failed to recover as the long range gradients present in the HuMoR TestOps method weren't present, and we found when blending rather than completely copying over that the weighting of the blending was another parameter that was difficult to choose, hence having a loss was the path of least resistance.

Next we experimented with jointly optimising x and z with a larger variety of losses. It was noted that the sitting down motion was not achieved consistently. However as previously mentioned, the initial z 's encode a sitting motion, thus we conject that the x 's and z 's were fighting and this sitting motion was lost, that is to say the z 's move to match the x 's motion in which the skeleton simply clips into the floor without bending legs, thus the z 's begin encoding less leg bending. To avoid this issue in future experiments, the z 's were detached from the consistency

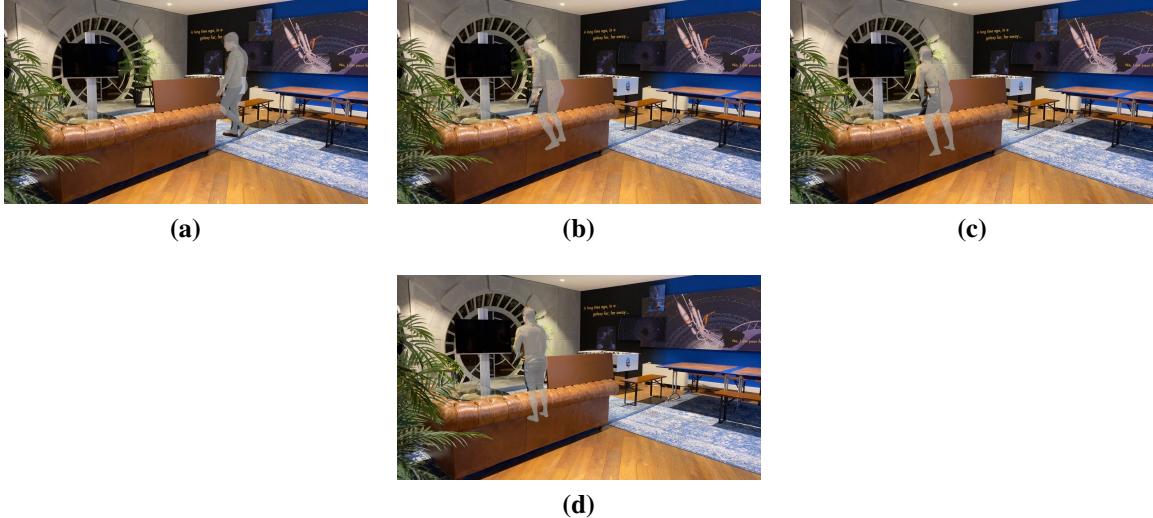


Figure 3.10: Stage 2 x 's

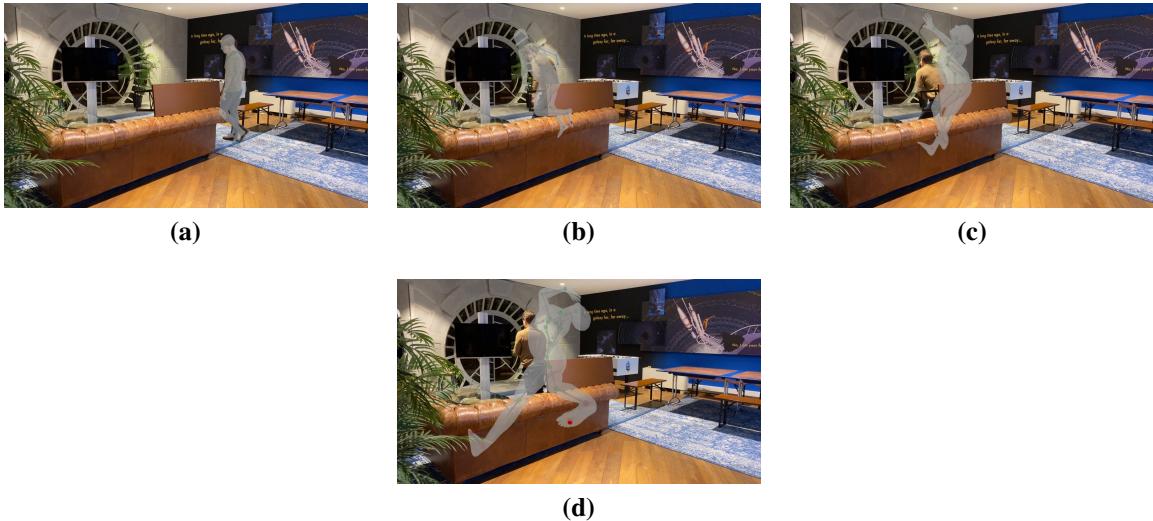


Figure 3.11: Deviation due to rollout of Stage 2 z 's

loss making the optimised x 's match the decoded x 's.

This led us on to some experiments in which we either optimise z whilst fixing x , or vice versa. With the z 's fixed at the stage 2 results, we managed to achieve a sensible sitting motion for a clip containing occluded sitting, but when we applied the same optimiser settings to a longer clip the results were no longer satisfactory. We also noted that when it did work, it required a large number of iterations, in the 1000s, which took 10mins or more, thus the speedup was not as evident as hoped (though our method would continue scaling better). In the inverse case, fixing x to the final HuMoR optimised states and optimising z , we noted that we consistently converge in the direction of the final HuMoR z 's, which indicates that we are optimising in the right direction, with the best results achieved at a rollout of 5 decoding steps.

All these experiments were performed with varying numbers of decoding steps, 1, 2, 5, and 10.

We found that the effect of more stable results but slower compute times with longer rollouts did manifest as intuition would suggest.

These experiments suggest that the next direction to pursue would be an optimisation scheme in which we flip flop optimising z then x . However considering our difficulty in achieving consistent results for the optimisation of the x 's with fixed z 's, our feeling that this might not be the optimal solution to the problem, my desire to try a new direction, and the time constraints, we deemed it time to move on and try a new method.

3.2.4 Conclusion

We found that jointly optimising the x 's and z 's did not prove an easy task to solve, which may indicate that

- it is a particularly difficult problem, and therefore that the optimiser struggles to satisfy the constraints and find a good minimum for both x and z at the same time
- the coupling of the x 's and z 's in the HuMoR decoder is an important issue. We found it to be at least a minor issue as it was necessary to decouple the z 's from the consistency loss (matching the optimised xs to the decoded x' s) to avoid the z 's being modified in such a way as to fit the unoptimised x 's at the beginning of the optimisation
- the optimiser is badly balanced

This at the very least suggests that a new decoder that does not couple x 's and z 's would be useful and worth investigating if a similar method is pursued.

In the investigations we found it possible to achieve sensible results on a given clip (for example fixing the unoptimised initial z 's, we can a plausible sitting motion in a small video of sitting in occlusion) but that the same settings for the optimiser failed to achieve a good result on an alternative clip. This may indicate that

- again, this is a difficult problem that the optimiser will struggle to solve
- the optimiser is badly balanced

as before.

While many of the issues encountered may simply be due to a badly balanced optimiser, considerable effort was made to mitigate this potential issue, and therefore it is our impression that this problem is at the very least a difficult one. We have improved our intuition, fixed several issues, and understood more about the problem and our goals. We therefore can make a more informed decision with respect to which direction to pursue next.

We found that taking a system that was designed with autoregression in mind, and trying to parallelise it, proved more difficult than expected. The experiments lead us to believe that jointly optimising a sequence of poses and latent variables obtained through a single frame decoder may be a more difficult formulation of the problem than others. We therefore conclude that a method in which we model longer motion sequences may be more fruitful.

TODO: Double check through Notes/Current_experiments.md to see if there's anything else we can learn from the experiments

Motion Diffusion Models

In Chapter ??., we came to the conclusion that the next sensible step would be an investigation of a model that operates directly on motion sequences. We also desired to have a general model that can be used for many different tasks, such that we might get a better overview of the uses of motion modelling. We therefore decided, due to promising literature presented in Section ??, that diffusion models were worth investigating.

4.1 Method

4.1.1 Diffusion Models

The diffusion framework that we use is defined theoretically as follows [?, ?, ?]. Given $x_0 \sim p(x_0)$ some data distribution, diffusion is a Markov noising process over latent variables $x_t, t = 1, \dots, T$, in which the process of adding noise (the 'forward noising process') is defined through the distributions

$$\begin{aligned} q(x_t|x_{t-1}) &:= \mathcal{N}(\sqrt{1 - \beta_t}x_{t-1}, \beta_t \mathbf{I}) \\ q(x_T|x_0) &:= \prod_{t=1}^T q(x_t|x_{t-1}) \end{aligned} \tag{4.1}$$

where with a sensible β schedule **TODO: Describe the beta schedule better** and a large T , we approximate a standard normal distribution as we approach $t = T$. This entails that the first latent variable x_T is approximately gaussian noise. The 'reverse noising process' is non tractable without the full data distribution, which we don't have access to, hence we approximate it using some parametrized function

$$p_\theta(x_{t-1}|x_t) = \mathcal{N}(\mu_\theta(x_t, t), \Sigma_\theta(x_t, t)) \tag{4.2}$$

4 Motion Diffusion Models

We find with p and q that this system forms a VAE that can be trained through the Variational Lower Bound (VLB) [?].

TODO: Put the classic figure of the noising/denoising process from [?] in here

In practice however we note some modification to this definition to ease our implementation and to improve performance. Firstly we note that the forward noising process can be written as

$$q(x_t|x_0) \sim \mathcal{N}(\sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)\mathbf{I}) \quad (4.3)$$

with $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$ and $\alpha_t = 1 - \beta_t$. This is simply a mixing of the data with gaussian noise according to a monotonically decreasing schedule $\bar{\alpha}_t$ and allows us to sample any latent at any timestep directly from the clean signal, which will aid our training procedure. We note that we take the following definition for the $\bar{\alpha}_t$ schedule, as in [?],

$$\begin{aligned} \bar{\alpha}_t &= \frac{f(t)}{f(0)} \\ f(t) &= \cos\left(\frac{t/T + s}{1+s} \cdot \frac{\pi}{2}\right) \end{aligned} \quad (4.4)$$

with the a small offset s (we take $s = 0.008$ as in [?]) to avoid β_t from being too small near $t = 0$ which can reduce performance according to the authors.

We can also present a more simple reverse noising process, where rather than predicting the mean and variance of the noise for a single step, we directly try to predict the denoised signal [?],

$$\hat{x}_\theta(x_t, t) \approx x_0 \quad (4.5)$$

which allows us to describe the denoising process as follows

$$p_\theta(x_{t-1}|x_t) = q(x_{t-1}|\hat{x}_\theta(x_t, t)) \quad (4.6)$$

TODO: check this definition of the denoising process against the code, I feel something is off
This direct prediction of the denoising signal allows us to use a the simplified MSE objective [?, ?] during training

$$\mathcal{L}_{simple} = \mathbb{E}_{x,t} [\|x - \hat{x}_\theta(x_t, t)\|_2^2] \quad (4.7)$$

which can still be interpreted as a form of Variational Lower Bound without the variance terms.

TODO: missing some note about the fixed small variance from the code?

4.1.2 Disney Motion Diffusion (DMD)

Armed with this understanding of diffusion, we define our Disney Motion Diffusion (DMD) model, strongly inspired by the work of Tevet. et. al [?].

State Representation

Our state represents a motion sequence of length N , and includes primarily a root translation, $r \in \mathbb{R}^{N \times 3}$, and 6d [?] joint orientations in body frame, $o \in \mathbb{R}^{N \times 6}$. The state can be extended

with some optional extras depending on the data, of foot ground contacts $b \in \mathbb{R}^{N \times 1}$, and explicit velocities for the root translation and joint orientations. Again, depending on the data, we can represent the root translation, and root orientations (a subset of the joint orientations) either in world space, where each motion sequence will have a distinct starting point and starting orientation, or as what we call a 'stuck root' reference frame, where the first pose is set to the positions $(0, 0)$ and identity orientation and all subsequent poses are described with respect to this first frame, with the view being that this formulation leads to a more general and smaller space that the model needs to understand.

Architecture

We follow closely [?] in our architecture, though without the text conditioning. We use a simple transformer [?] encoder as the denoising network, conditioned on the diffusion timestep so that the network can learn how much noise is present that needs removing. This is used with the as described denoising process which can be seen in Figure ??.

TODO: Link to appendix with more details on the architecture choices, number of hidden layers, etc

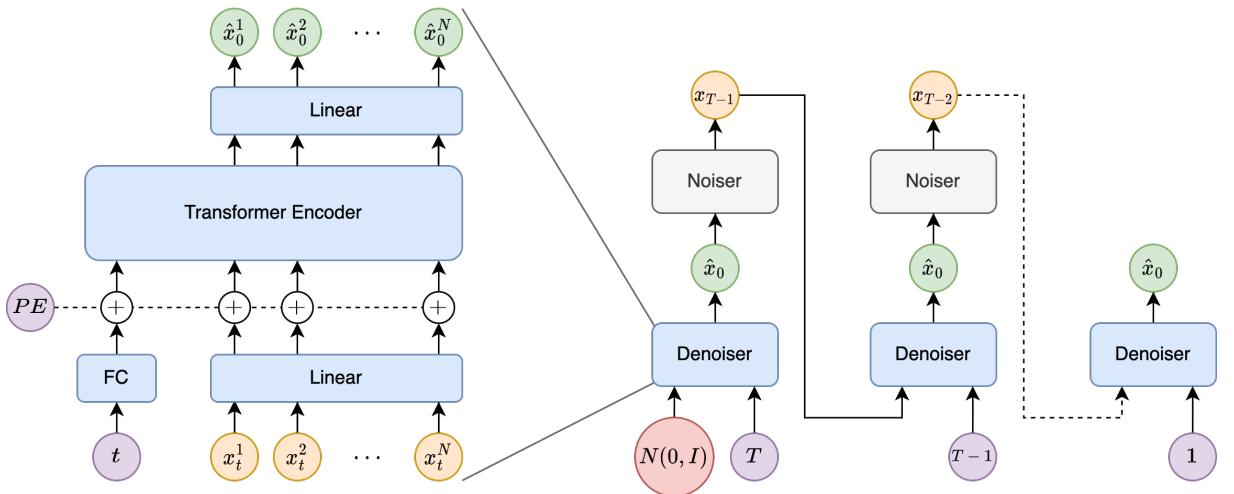


Figure 4.1: DMD Architecture & Denoising Process

Training

In addition to the simple MSE loss described in section ??, we can include a number of other losses to aid the training procedure.

$$\begin{aligned}\mathcal{L}_{fk} &= \frac{1}{N} \sum_{i=1}^N \|x_0^{(pos)} - FK(\hat{x}_0^{(i)})\|_2^2 \\ \mathcal{L}_{vel} &= \frac{1}{N-1} \sum_{i=1}^{N-1} \|(x_0^{(i+1)} - x_0^{(i)}) - (\hat{x}_0^{(i+1)} - \hat{x}_0^{(i)})\|_2^2 \\ \mathcal{L}_{contact} &= \frac{1}{N-1} \sum_{i=1}^{N-1} \left\| \left(FK(\hat{x}_0^{(i+1)}) - FK(\hat{x}_0^{(i)}) \right) \cdot \hat{b}^{(i)} \right\|_2^2\end{aligned}\quad (4.8)$$

with $x_0^{(i)}$ representing the i th pose of the motion sequence x_0 . The \mathcal{L}_{fk} loss is a forward kinematics loss, i.e that where the positions are calculated from the orientations, and compared to the ground truth positions that we have access to, $x_0^{(pos)}$, and which provides another mechanism by which the orientations might be adjusted, thus increasing the robustness of the predictions. The \mathcal{L}_{vel} loss is a calculated velocity loss, in which the first order finite differences are used to estimate the velocities of the state (orientations and root translation only in this case), and which encourages the model to learn a smooth motion sequence. The $\mathcal{L}_{contact}$ loss is a foot contact loss, where the velocities of the ground contacting feet joints, indicated by the predicted foot contacts $\hat{b}^{(i)}$, are minimised so as to discourage foot sliding. We note that the $\mathcal{L}_{contact}$ loss is only used when the ground contacts are included in the state representation.

Inpainting

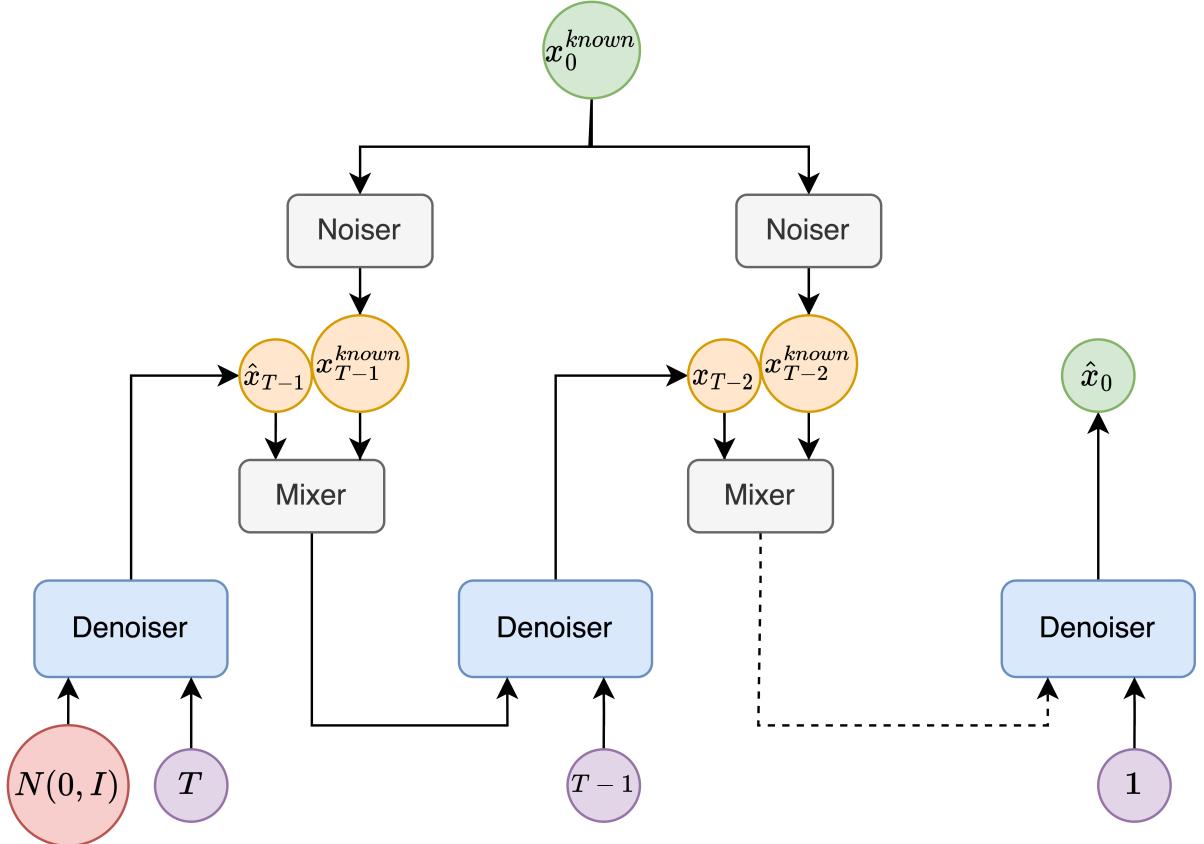
An important tool in the arsenal of the diffusion model literature is that of inpainting [?]. Inpainting is the process of keeping a portion of the state constant throughout the diffusion process. An example of the use of this technique would be that of changing the motion of the lower body while maintaining the motion of the upper body, though there are many other use cases that are explored in more detail in Section ???. The inpainting procedure works as follows; at each denoising step we mix a portion of the latent variable state \hat{x}_t , indicated through the binary mask m , with the complementary portion of some known state $x^{(known)}$, indicated through $1 - m$,

$$\hat{x}_t = m \odot \hat{x}_t + (1 - m) \odot q(x^{(known)} | t - 1) \quad (4.9)$$

Graphically this process can be seen in Figure ???. This procedure is particularly attractive at it requires no special training procedure, and can be employed readily with any diffusion model.

Data

At Disney Research|Studios we have access to professionally capture motion capture data that encompasses a wide range of possible motion for a number of subjects. In total we have TODO: X motion sequences, varying in length from TODO: X to TODO: X seconds long, thus in total

**Figure 4.2: Inpainting Procedure**

we have **TODO: X** minutes of motion capture data. This data is normalised elementwise and retargeted to a standard skeleton. A number of motion sequences are held out for testing.

TODO: Describe dataset 1 (fixed root, not contacts, etc) and 2 (moving root, stuck root, contacts, how to get contacts, etc.)

4.2 Experiments

Our main goal in this investigation into diffusion models was to get a better understanding of their performance on a wide range of tasks. We keep in mind the original task of improving an existing motion sequence, but are aware of the other potentials of the model and so do not limit ourselves in scope.

4.2.1 Test Cases

We present here a variety of different test cases that on which the models performance can be evaluated. The general structure of the evaluation is to generate a sufficiently large number of sequences for each given task, then to import them into blender and visually inspect their quality and how well they achieve each task. We have decided to take a qualitative approach

4 Motion Diffusion Models

as we are attempting to get a better understanding of the models capabilities, and so this was the path of least resistance to obtain such an understanding in the time that was left after the HuMoR investigation. We note that in the future, if a given task is deemed worth investigation further, or if a significantly larger number of models are being tested, more rigorous evaluation metrics would be greatly beneficial.

Sequence Generation

To ensure that the model has learned a sensible variety of human motion, we generate a number of sequences from gaussian noise, and check that they contain sufficient diversity, considering that a main indicator of an overfit diffusion model is that of a lack of diversity.

Denoising

TODO: Come up with a better test for this

Occluded Legs

An important situation in which the current Disney Research|Studios model does not perform optimally is that of occluded sequences. To test the ability of the diffusion models to handle this situation, we evaluate the models capabilities to generate the leg motion of an entire motion sequence. We achieve this using the inpainting framework described in Section ??, keeping the upper body joint rotations and the root translation fixed, and allowing the model free reign over the rest of the state.

Missing State

To further test the ability of the models to handle occlusions, a test case where we only keep a random percent of the state during inpainting. This can either be

1. a random subset of each frame in a motion sequence
2. a random subset that is constant across the whole sequence

with the first testing the models ability to handle arbitrarily occurring missing data, and the second testing the models ability to recover from longer term occlusions.

Inbetweening

Another important area of motion generation research is that of Motion Inbetweening, in which the goal is to fill in a middle part of a motion sequence in which a number of frames at the beginning and end of the sequence are given. Again we can use the inpainting framework from Section ?? and keep the edges of the motion as static.

4.2.2 Baseline Model

The baseline model was trained with the Disney Research!Studios dataset that does not, and cannot (due to a lossy rendering of the mocap), contain foot contacts. The state includes the root translation, joint orientations, and their respective velocities explicitly. It was trained for **TODO: X** epochs. **TODO: More details?**

Sequence Generation

We found that the model generates a diverse set of plausible motion, a number of frames are shown in Figure ???. A few shortcomings however can be seen. Most notably we find that the model produces a fair bit of foot sliding **TODO: Is this true??**, where the foot seem to be in contact with the ground but has a non zero velocity in some directions. This effect can be mitigated through the use of the foot sliding loss coupled with data with containing labeled foot contacts.

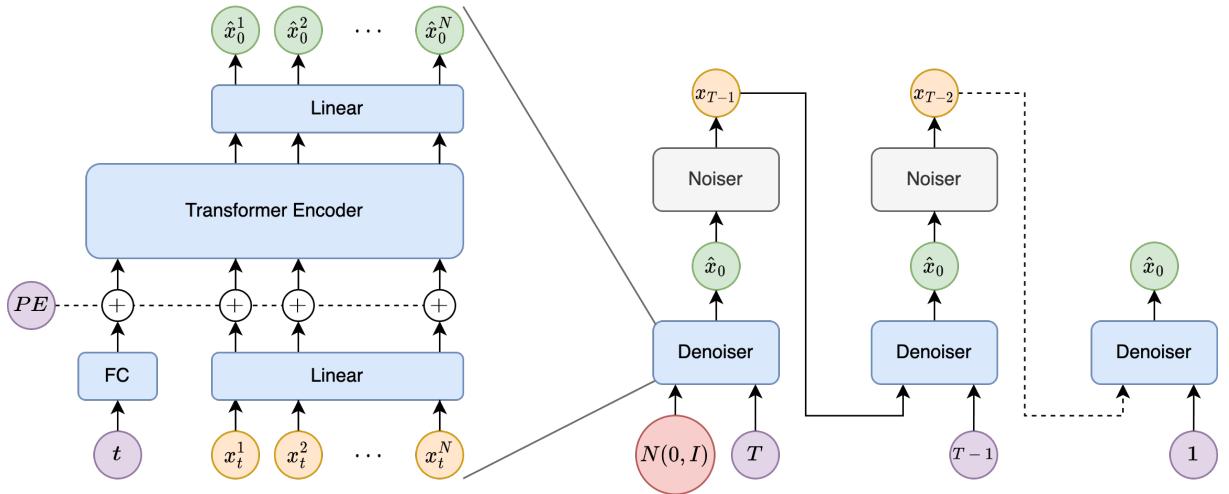


Figure 4.3: Motion Generation - Baseline model

Denoising

TODO: Come up with a better test for this

Occluded Legs

The model shows a strong performance on the task of inpainting missing legs in a motion sequence, with the result being of leg motion that is often both smooth and well adhering to the upper body motion. A few examples frames are shown in Figure ???. The most notable failure case is that sometimes the model reverses the direction of motion. We postulate that this is due to the fact that the data for the baseline model has a root that does not move in the forward/reverse direction, therefore the pendulating rotation motion of the hips, without the additional context of a specific direction, can often be both interpreted as both forward and backward motion.

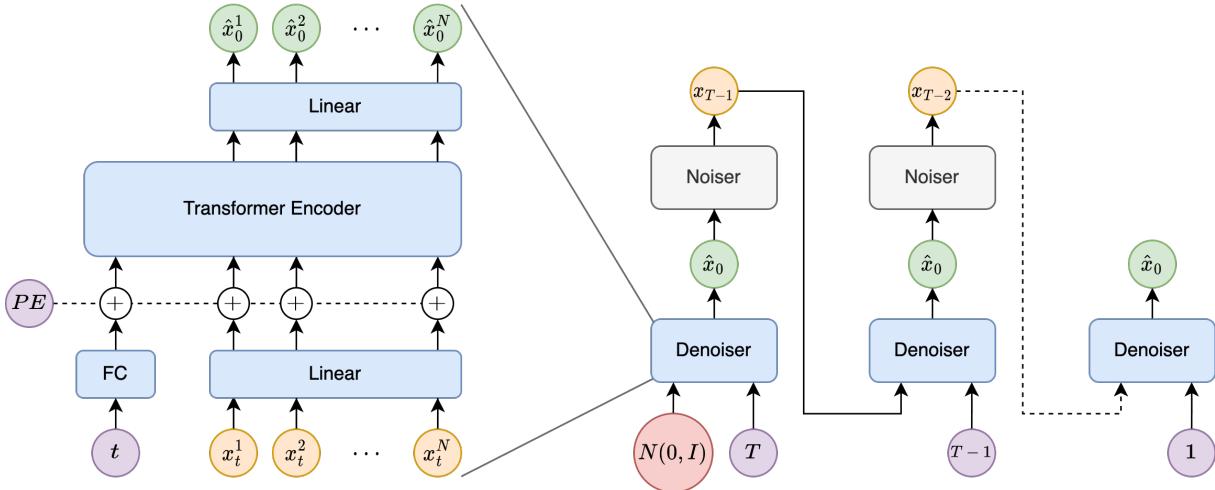


Figure 4.4: Occluded Legs Inpainting - Baseline model

Missing State

Random each frame When removing a new random subset of the joints of each pose every frame in the motion sequence, we find that the model can very accurately reconstruct the missing state, even when only 50% of the state is given each frame, and surprisingly it still performs reasonably well when only 25% of the state is given each frame. A few example frames are shown in Figure ???. This shows that the model has learned well the concept of temporal consistency, as in this experiment it is likely that a missing joint will be surrounded (temporally) by given joints, hence it simply needs to interpolate between the given joints.

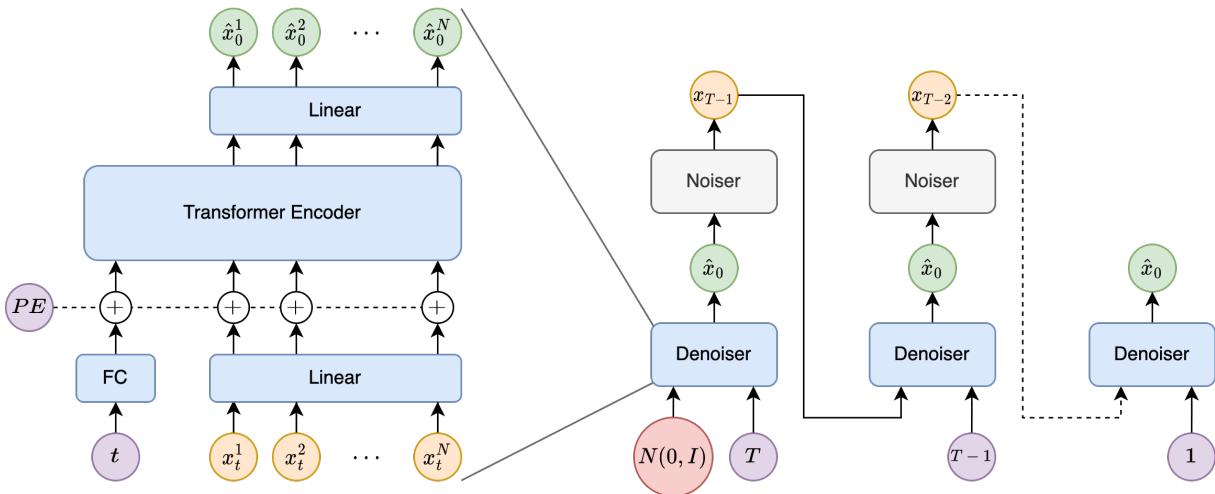


Figure 4.5: Occluded Legs Inpainting - Baseline model

Random for the whole sequence **TODO: Describe better**

Inbetweening**4.2.3 Model with Foot Contacts****Sequence Generation****Denoising****Occluded Legs****Missing State****Inbetweening****4.2.4 Autocompletion**

Another project that is running at Disney Research|Studios is that of pose autocompletion. This is the task of taking a set of *handles*, a subset of joints, and predicting the position of the rest of the joints. The goal here is to again help to improve the animation pipeline, by reducing the amount of work an animator has to do to get a sensible pose. As a side note to the previous experiments, we trained some diffusion models for the described task, where we enforce a sequence length of 1 so that we are only predicting single pose, and where we modify the inpainting to allow for a specific set of joints to be deemed *handles* and to be kept fixed.

Baseline

The model is as described in Section ?? however with a sequence length of 1.

We find that while the model does present plausible poses, and that on average these poses have handle joints close to the desired positions of the handles (with a few exceptions), it does not satisfy the task properly as can be seen in **TODO: include images of position of handles vs. position of actual joints**. An ideal solution would not move the handles at all, and would simply modify the rest of the joints.

TODO: include images

Handle Conditioning

To mitigate this issue of the handles being modified, we propose a handles conditioning, in which a one hot embedding of the handles is run through a linear layer, then concatenated to the input of the network. A loss is then applied to the output that penalises any deviation in position of the handles that are output as compared to those input. In theory this should force the network to learn to consider the handles as immutable, and to create a pose that satisfies the handles.

Though the implementation was largely completed, we unfortunately did not have time to finish debugging some of the finer details, and so cannot provide concrete results. However we deemed the idea interesting enough to present in this thesis, and hope that it can be of use in future research directions.

4.3 Conclusion

The DMD diffusion framework described in Section ?? proves to represent an exciting research direction. The baseline model shows good performance on a variety of tasks, as seen in Section ??, demonstrating that the diffusion framework is able to learn a general purpose motion prior model without task specific fine tuning or architecture design. This gives us hope for future research directions in which the specific task to solve is more precisely described, where more fine grained architecture, data and training choices could be made that could further improve the performance of the models.

The implementation of the diffusion framework within the Disney Research|Studios codebase paves the way for further investigation into its uses in the context of general animation tasks.

TODO: Longer conclusion?

4.4 Future Work

4.4.1 Model

Some of the shortcomings of the model, seen in some of the Inpainting tasks and most notably in the Inbetweening task, is that of inconsistency between the specified known motion and the motion subsequently generated during inpainting by the model. For inbetweening, if we avoid overwriting the generated motion with the known motion in the last denoising step, the motion is very smooth **TODO: image**, however does not quite match the known motion, as when we override the generated motion with the known motion, we find that there is a discontinuity between the generated and known motion, i.e the motion jumps. This effect might be mitigated with a number of strategies. Firstly we might condition the model with an indication of where a known area is specified and penalise it during training for deviating from this known motion, as described in Section ???. Another method would simply be that of a postprocessing step, in which the known motion and the generated motion is mixed for a number of frames, thus achieving a smooth transition between the two.

The model training could also be improved. The stopping criteria is currently taken as the point where the loss visually seems to have converged on the loss graphs. This however is not a robust method as it does not take into account the possibility of overfitting or the possibility of decreasing performance. This can be mitigated through the use of validation metrics as described in Section ???, or by a method such as that presented in [?], in which they retrospectively choose the checkpoint that minimizes the FID score.

4.4.2 Model use

- Inpainting - We notice that the target changes often, even in the last steps - Maybe change the schedule to repeat the last steps a number of times to allow more time for the process to converge - Autocomplete

4.4.3 Model evaluation

- We presented an entirely qualitative evaluation of the models - We should make more effort to report on the metrics - Metrics - Generation - FID score? (the one comparing the distribution to the a realistic dist of data) - Diversity - Checking it's not just regurgitating the training data
- Inbetweening - Smoothness around the borders - Denoising/missing joints - Difference to the ground truth - [?] - Physical foot contact score - Model comparison - We should compare the models to some baseline models - For example the missing state evaluation where a new subset of the joints is removed each frame might well be accurately reconstructed by simple linear interpolation between the missing frames, so the models performance might not be so impressive afterall.

4.4.4 Data

TODO: Cite some paper that says data is more important than anything We note that our dataset, with **TODO: X** minutes of motion spanning **TODO: X** subjects is significantly smaller than the AMASS [?] dataset that contains more than 40hours of motion data spanning 300 subjects. We therefore expect that our model will generalise less well for than if it were trained on the AMASS dataset. We also note however that this is not possible due to the licensing restrictions on the AMASS dataset limiting it to non-commercial use, and so if any of the models we explored in this thesis or in subsequent follow up work are to be productionized, it would be prudent to obtain a larger dataset of motion, especially considering that the use cases in which such a system would be most helpful to an animator might well be those that are furthest from the data distribution we currently have, as the more obscure and intricate the motion sequence the harder it would be to animate.

Conclusion and Outlook

TODO: Chpt 1

Appendix

Appendix intro

A.1 Section

TODO: