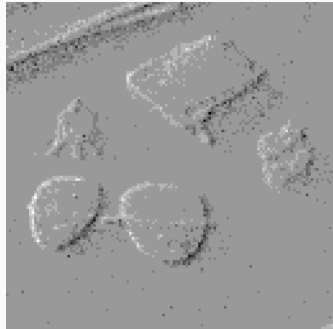


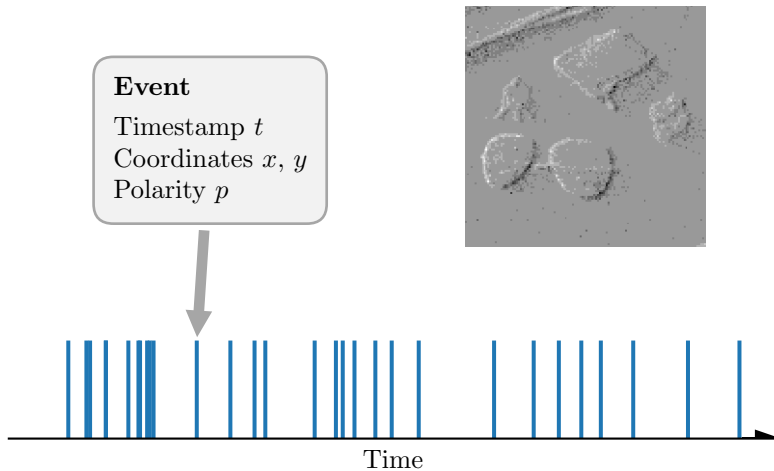
Gesture Recognition with a Neuromorphic Vision Sensor and Deep Learning

Marten Lienen

Image vs. Event Histogram



The Event Stream



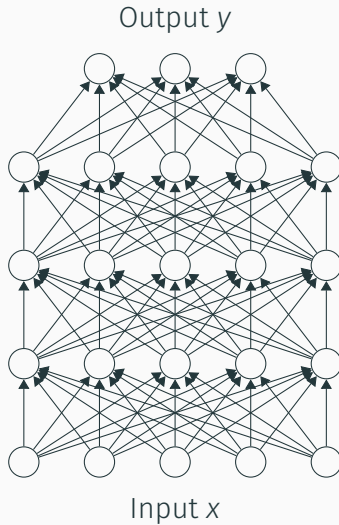
Dataset

- 16 gestures
- 2 subjects
- 40 segmented recordings
- 40 minutes and 640 examples in total

Preprocessing

- Translation invariance in time $\Delta t_i \leftarrow t_i - t_{i-1}$
- Translation invariance in space
 - Means with exponentially decaying weights as point of reference
 - Halflife of 50 ms and 1 s
- Feature Normalization

Neural Networks



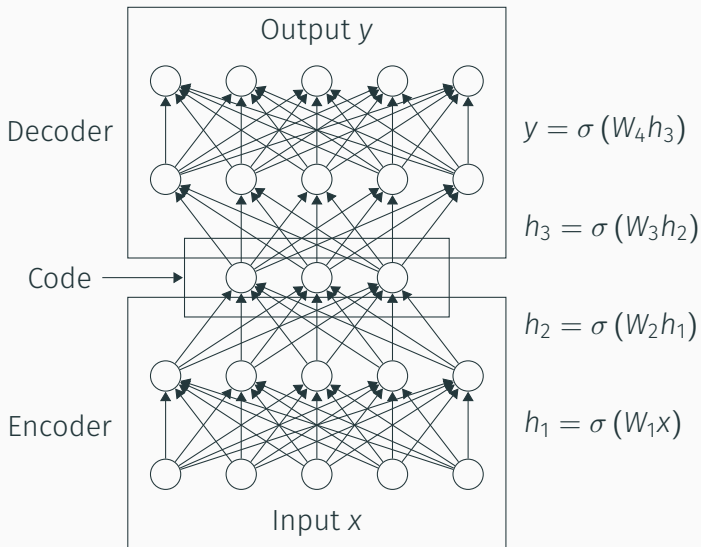
$$y = \sigma(W_4 h_3)$$

$$h_3 = \sigma(W_3 h_2)$$

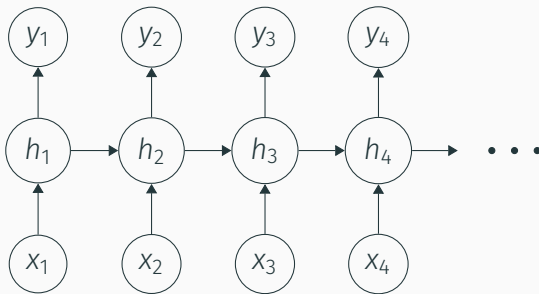
$$h_2 = \sigma(W_2 h_1)$$

$$h_1 = \sigma(W_1 x)$$

Autoencoders

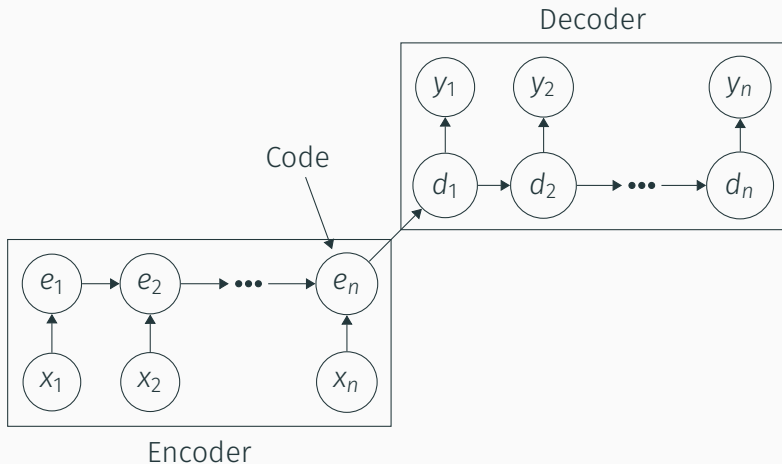


Recurrent Neural Networks

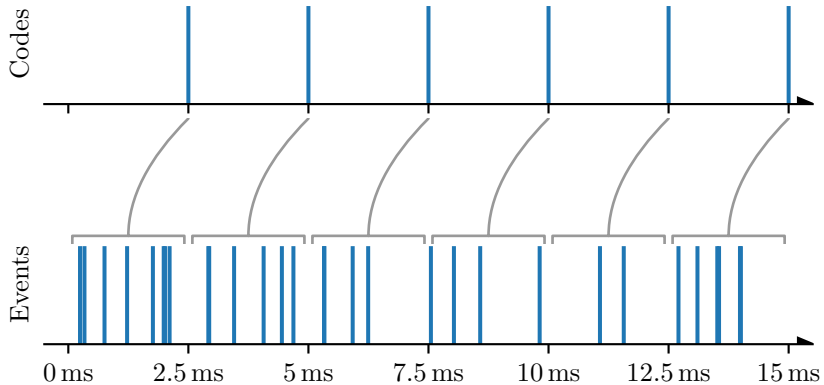


$$h_t = \sigma(Wx_t + Uh_{t-1})$$

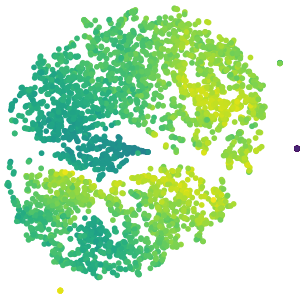
Recurrent Autoencoders



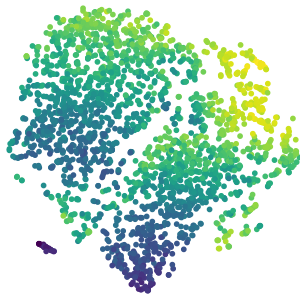
Event Sequence Compression



Learned Representations

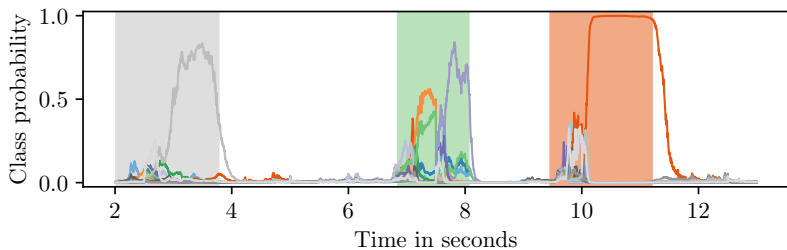


(a) Push Hand Up

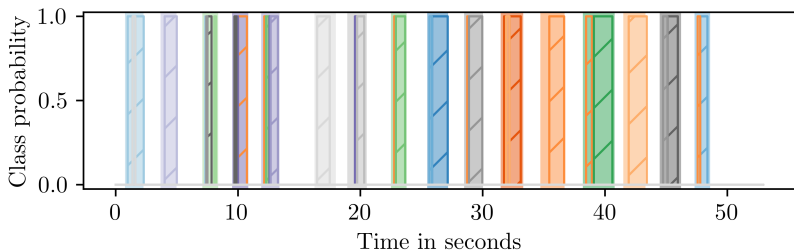
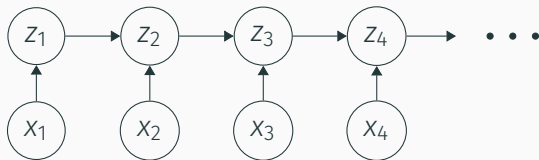


(b) Swipe Right

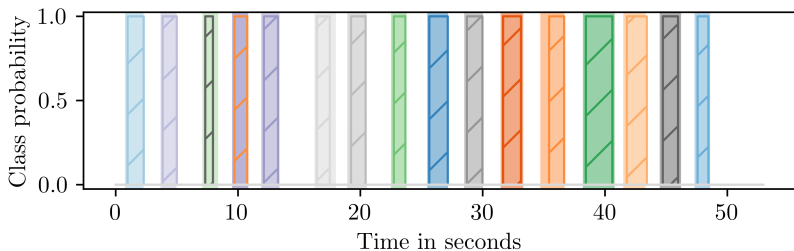
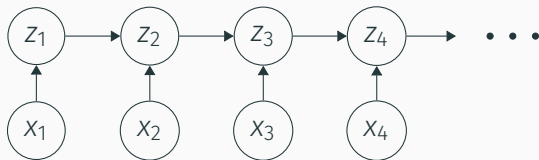
Framewise Classification



HMM Decoding



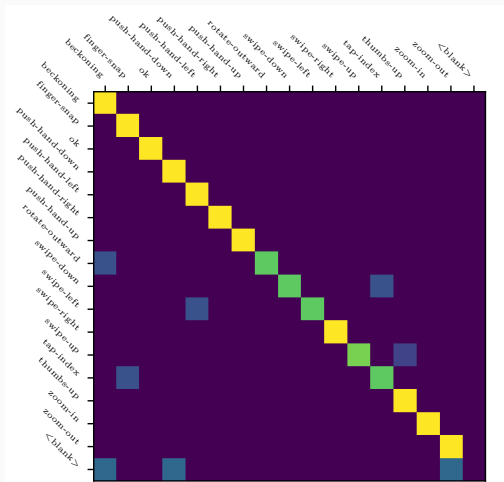
HMM Segmentation



Predictions

Predicted Sequence	True Sequence
finger-snap	rotate-outward
ok	ok
zoom-in	zoom-in
swipe-left	swipe-left
swipe-down	swipe-down
zoom-in ✗	zoom-out
zoom-out	beckoning
beckoning	push-hand-right
push-hand-right	swipe-right
swipe-right	push-hand-down
push-hand-down	push-hand-up
push-hand-up	thumbs-up
thumbs-up	swipe-up
swipe-up	tap-index
tap-index	finger-snap
tap-index	

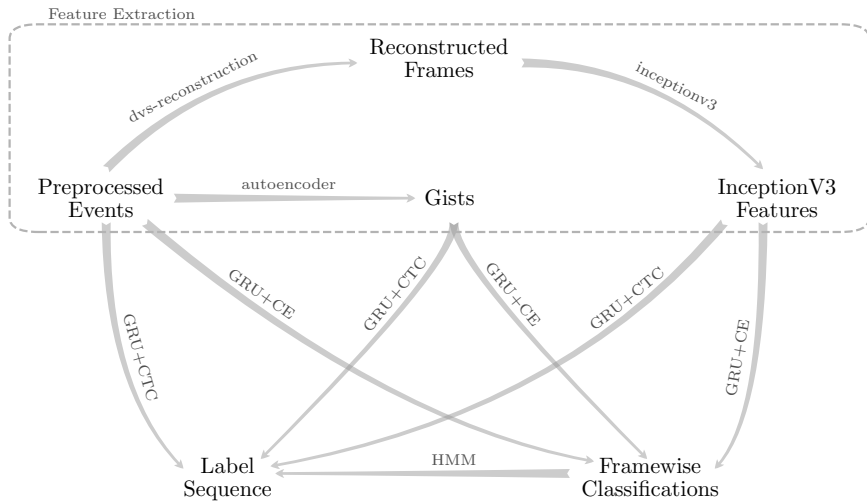
Results



Mean Levenshtein Distance 2.0, Label Error Rate 0.125

Thank You

Methods



Frame Reconstructions

