

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053

email

[illegible]

2 Prior Work

3 Offset Pooling

$$m_k = \int_{N_k} z(f, x, y) \frac{e^{\beta z(f, x, y)}}{\int_{N_k} e^{\beta z(f', x', y')} dv'} dv \approx \max_{N_k} z(f, x, y) \quad (1)$$

1

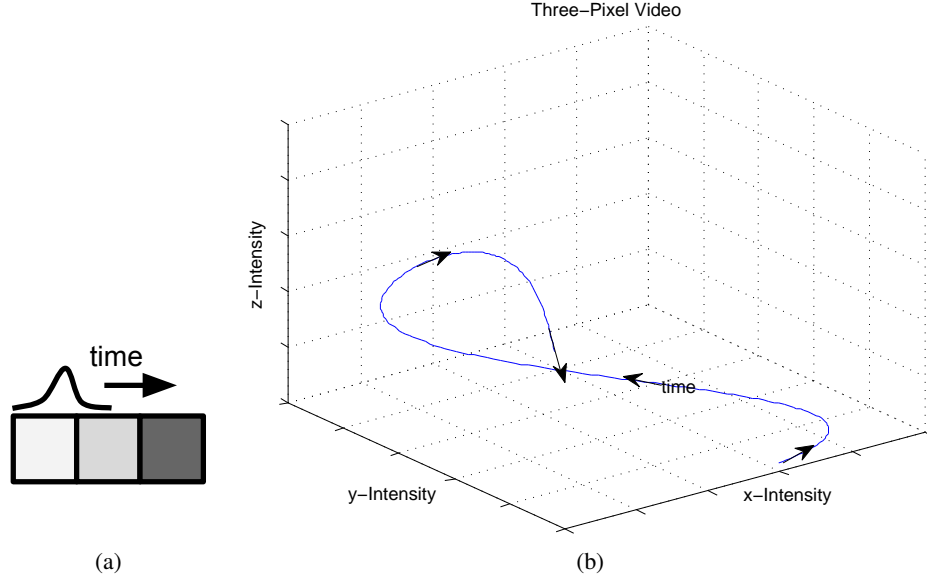


Figure 1: (a) A video generated by translating a Gaussian intensity bump over a three pixel array (x,y,z) , (b) the corresponding manifold parametrized by time in three dimensional space

three components. In general, the number of components in \mathbf{p}_k is equal to the dimension of the topology of our feature space induced by the pooling neighborhood. The dimensionality of \mathbf{p}_k also has the interpretation as being the *maximal intrinsic dimension of the data*. If we define a local standard coordinate system in each pooling volume which is bounded between -1 and +1, the soft “*argmax - pooling*” operator is defined by the vector-valued integral:

$$\mathbf{p}_k = \int_{-1}^1 \begin{bmatrix} f \\ x \\ y \end{bmatrix} \frac{e^{\beta z(f,x,y)}}{\int_{N_k} e^{\beta z(f',x',y')} dv'} dv \approx \arg \max_{N_k} z(f,x,y) \quad (2)$$

To motivate the definition of these operations consider a video generated by translating a Gaussian “intensity bump” over a three pixel array at constant speed. The video corresponds to a one dimensional manifold in three dimensional space, i.e. a curve parameterized by time. Next, assume that some feature detector activates when centered on such a bump. Applying the *max*-pooling operator over this region reveals that the Gaussian bump is present in somewhere in this region (i.e. *the what*). Applying the *argmax* pooling operator over the region returns the position (i.e. *the where*) with respect to a local coordinate frame defined over the pooling region. This position variable varies linearly with respect to time, and locally parameterizes the data manifold. More generally offset-pooling can potentially locally parametrize arbitrary transformations if the features we pool over are learned. Moreover, linear predictions can be formulated in the parameter space provided that there is a mapping from the parameter space back to the input space (decoder).