# Recurrent Neural Networks for Classifying Human Embryonic Stem Cell-Derived Cardiomyocytes

## Carolina Pacheco and René Vidal

Center for Imaging Science, Mathematical Institute for Data Science, Department of Biomedical Engineering, Johns Hopkins University, Baltimore, USA
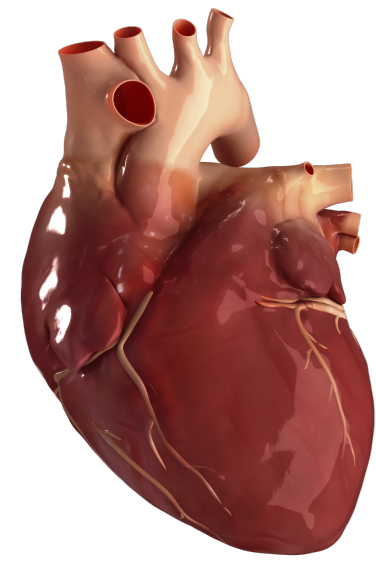
## Motivations

- Classifying human embryonic stem cell-derived cardiomyocytes (hESC-CMs) is important in cardiac regenerative medicine to reduce the risk of complications with stem cell therapy.
- Adult CMs can be classified based on the characteristic shape of their action potentials (APs), but the understanding of how the shape of hESC-CM APs relates to that of adult CMs remains limited [4].
- Prior approaches to hESC-CM APs classification either discard most of the information contained in the APs or rely on distances that are computationally expensive [5].
- Our goal is to develop a classifier of hESC-CM APs applicable to large-scale datasets.
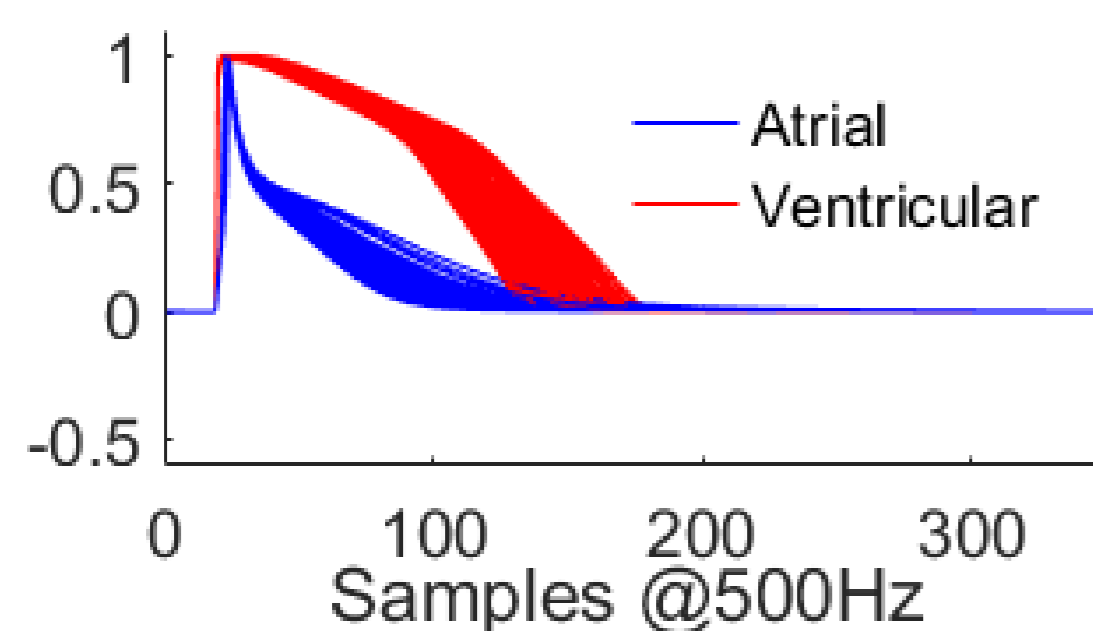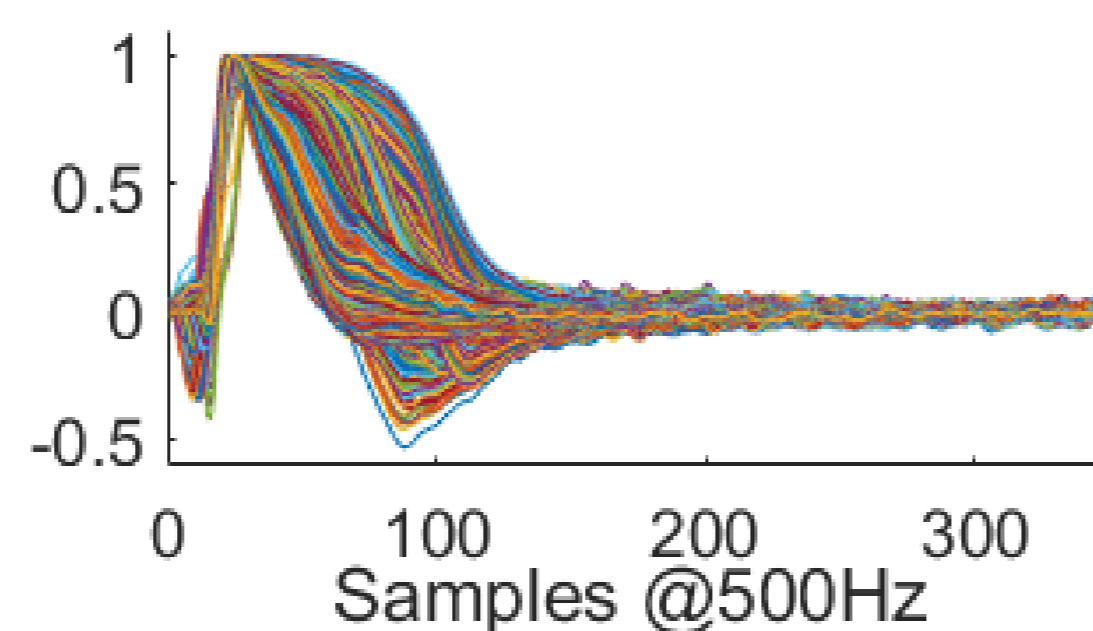


Fig. 1: APs from adult CMs



Fig. 2: APs from hESC-CMs

## Contributions

- We propose a new method for classifying hESC-CM APs based on recurrent neural networks (RNNs) with long short term memory (LSTM) units.
- The learning approach exploits the abundance of labels for adult APs, which can be obtained via simulation of electrophysiological models of the typical adult phenotypes.
- Our semi-supervised approach uses a novel loss function that combines a supervised classification loss for adult APs and an unsupervised contrastive loss for hESC-CM APs.

## Problem Formulation

- Set of unlabeled hESC-CM APs $\Omega_e = \{\mathbf{x}_j^e\}_{j=1}^{N_e}$: each $\mathbf{x}_j^e$ is a time-series of length $K$.
- Set of labeled adult CM APs $\Omega_a = \{\mathbf{x}_i^a, y_i^a\}_{i=1}^{N_a}$: each $\mathbf{x}_i^a$ is a time-series of length $K$ labeled as atrial ($y = 0$) or ventricular ($y = 1$).
- Problem: Assign a label $\hat{y}^e$ to a new $\mathbf{x}^e$, where $\hat{y}^e = 0$ denotes atrial-like and $\hat{y}^e = 1$ denotes ventricular-like.

## Classifier Architecture

- RNN with LSTM units as a classifier
  - Hidden layer:
    LSTM of dimension $p = 3$
  - Output layer:
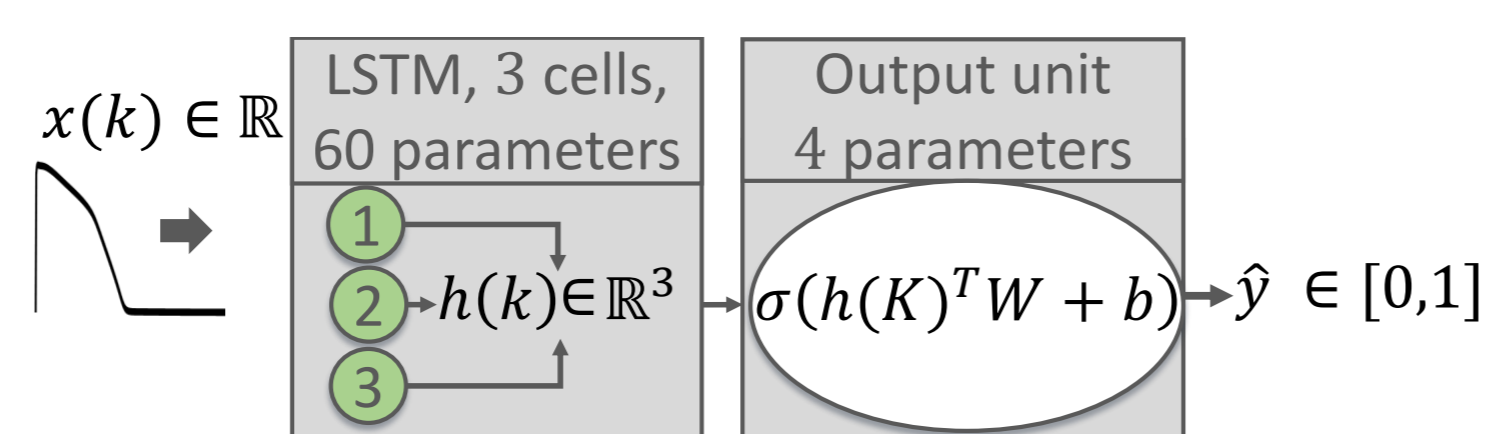    Sigmoid unit $\sigma(z) = \frac{1}{1+e^{-z}}$



Fig. 3: Network architecture

## Semi-Supervised Loss Function

The supervised part guides the LSTM to correctly predict labels of adult CMs, while the unsupervised part guides the LSTM to predict the same label for similar embryonic CMs:

$$\frac{1-\lambda}{N_a}\left(\sum_{j=1}^{N_a}\ell_s\left(y_j^a, \hat{y}_j^a\right)\right) + \frac{\lambda}{N_e-1}\sum_{j=2}^{N_e}\ell_u\left(\hat{y}_j^e, \hat{y}_{j-1}^e\right).$$

- Binary crossentropy loss on adult CM APs
$$\ell_s(y_i^a, \hat{y}_i^a) = -y_i^a \log(\hat{y}_i^a) - (1 - y_i^a)\log(1 - \hat{y}_i^a).$$

- Contrastive unsupervised loss on hESC-CM APs
$$\ell_u(\hat{y}_j^e, \hat{y}_{j'}^e) = s_{(j,j')} \cdot \ell_s(\hat{y}_j^e, \hat{y}_{j'}^e) + \left(1 - s_{(j,j')}\right) \cdot \ell_s((1 - \hat{y}_j^e), \hat{y}_{j'}^e),$$

where $s_{(j,j')} = \exp\left(-\frac{d^4(\mathbf{x}_j^e, \mathbf{x}_{j'}^e)}{\sigma_s^4}\right) \in [0,1]$ represents the similarity between $\mathbf{x}_j^e$ and $\mathbf{x}_{j'}^e$

- Euclidean distance $d(\mathbf{x}_j^e, \mathbf{x}_{j'}^e) = \frac{1}{\sigma_M}\sqrt{\sum_{k=1}^{K}\left(x_j^e(k) - x_{j'}^e(k)\right)^2}$, or
- Metamorphosis distance
$$d(\mathbf{x}_j^e, \mathbf{x}_{j'}^e) = \sqrt{\min_{\mathbf{x},\mathbf{v}}\sum_{r=0}^{R-1} \|v(k,r)\|_{V_d}^2 + \frac{1}{\sigma_M^2}\|x(k+v(k,r), r+1) - x(k,r)\|_2^2},$$

where $x(k,r)$ is the interpolation path between $x(k,0) = \mathbf{x}_j^e(k)$ and $x(k,R) = \mathbf{x}_{j'}^e(k)$, $\sigma_M$ and $\sigma_s$ are normalization parameters, $\mathbf{v}$ is the velocity of the interpolation path and $\|\cdot\|_{V_d}^2$ is a Sobolev norm.

## Clustering Quality Index

- Davies-Bouldin Index (DBI) [1]. The DBI between two clusters $\Omega_0 = \{\mathbf{x}_j^e \mid \hat{y}_j^e < 0.5\}$ and $\Omega_1 = \{\mathbf{x}_j^e \mid \hat{y}_j^e \geq 0.5\}$ is the ratio between intra-cluster dispersion and distance between clusters

$$DBI(\Omega_0, \Omega_1) = \frac{S_0 + S_1}{M_{01}},$$

where $S_y$ is the mean distance from elements of class $y$ to the average signal of the same class, and $M_{01}$ is the distance between the average signals of both classes.

## Experiments

- Dataset
  - Adult CM APs: 300 synthetic adult APs generated using the O'hara-Rudy ventricular model (ORd) [2] and the Nygren atrial model [3].
  - hESC-CM APs: 6940 unlabeled hESC-CM APs obtained from 9 cell aggregates [4].
- Implementation Details
  - Keras with TensorFlow backend, RMSProp optimizer, batches of 3 adult APs and 16 hESC-CM APs (90 batches validation, 10 batches training).
- Baselines
  - 1NN E. 1-Nearest-Neighbor method with Euclidean distances [5]
  - 1NN M. 1-Nearest-Neighbor method with metamorphosis distances [5]
  - 1NN E SMRS. 1-Nearest-Neighbor method with Euclidean distances using 300 templates
- Results

Three cases studied:
- Sup-LSTM. Supervised learning $\lambda = 0$
- Semi-LSTM-E. Semi-supervised learning $\lambda = 0.1$ with Euclidean distances
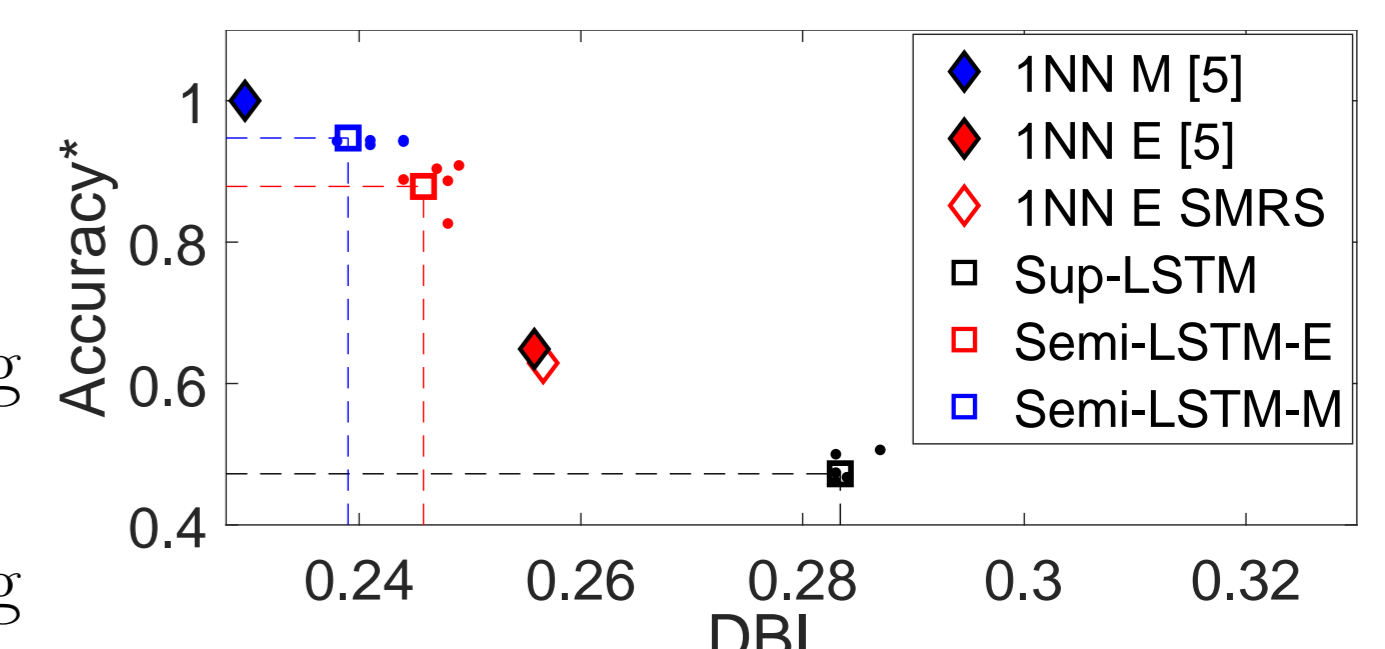- Semi-LSTM-M. Semi-supervised learning $\lambda = 0.1$ with metamorphosis distances



Fig. 4: Accuracy* vs DBI. 1NN M as ground truth (E: Euclidean, M: Metamorphosis).

Computational time for classifying a new sample:
- 1NN M. 6.74 sec/sample in 2 8-core computer nodes with 8 2.3GHz CPUs per node.
- Semi-LSTM. < 6 sec for the whole 6940 APs dataset in one 2.2 GHz CPU with 2 cores.
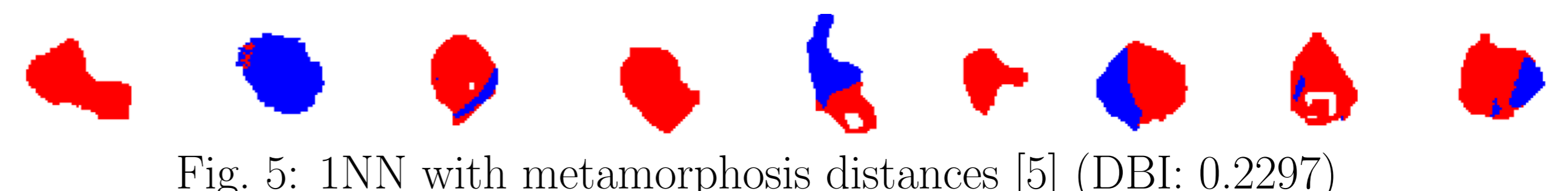


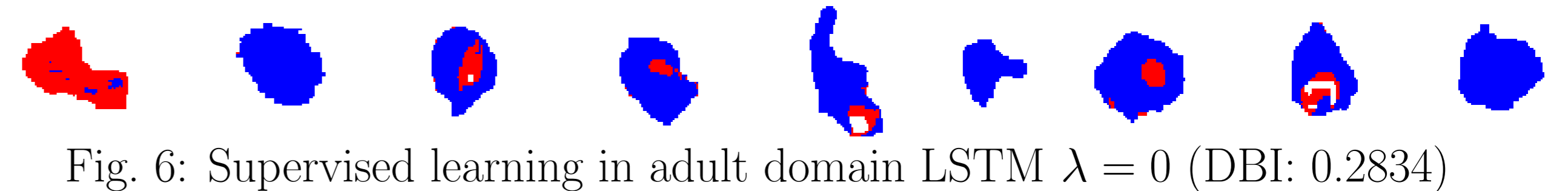Fig. 5: 1NN with metamorphosis distances [5] (DBI: 0.2297)



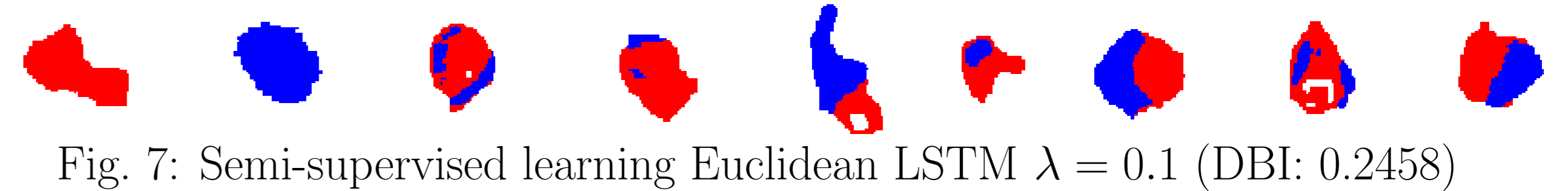Fig. 6: Supervised learning in adult domain LSTM $\lambda = 0$ (DBI: 0.2834)



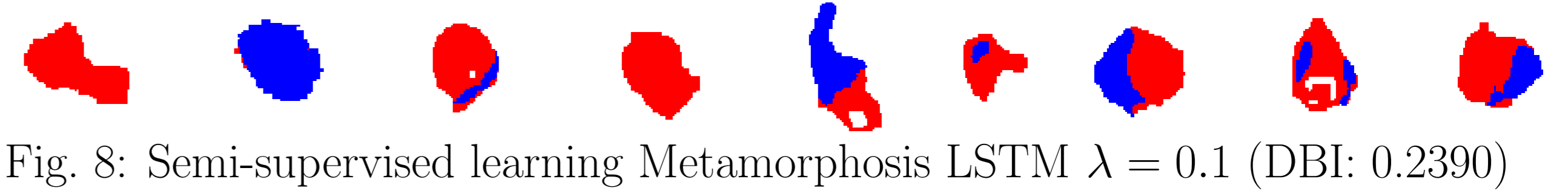Fig. 7: Semi-supervised learning Euclidean LSTM $\lambda = 0.1$ (DBI: 0.2458)



Fig. 8: Semi-supervised learning Metamorphosis LSTM $\lambda = 0.1$ (DBI: 0.2390)

## Conclusions

- Experiments confirm the benefits of integrating information from both adult and embryonic cardiomyocytes in a semi-supervised learning scheme for hESC-CMs classification.
- The proposed semi-supervised approach uses the Euclidean metric more effectively than previous methods, outperforming the 1NN scheme.
- Proposed semi-supervised approach gives results similar to the state-of-the-art (94.73% of agreement) with clear computational advantages when applied to new samples.

[1] D.L. Davies, and D.W. Bouldin. A cluster separation measure. IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 2, pp. 224-227, 1979.

[2] T. O'Hara, L. Virág, A. Varró, and Y. Rudy. Simulation of the undiseased human cardiac ventricular action potential: model formulation and experimental validation. PLoS Computational Biology, vol. 7, no. 5, p. e1002061, 2011.

[3] A. Nygren, C. Fiset, L. Firek, J. Clark, D. Lindblad, R. Clark, and W. Giles. Mathematical model of an adult human atrial cell: the role of k+ currents in repolarization. Circulation Research, vol. 82, no. 1, pp. 63-81, 1998.

[4] R. Zhu, M.A. Millrod, E.T. Zambidis, and L. Tung. Variability of action potentials within and among cardiac cell clusters derived from human embryonic stem cells. Scientific Reports, vol. 6, p.18544, 2016.

[5] G. Gorospe, R. Zhu, J-Q. He, L. Tung, L. Younes, and R. Vidal. Efficient metamorphosis computation for classifying embryonic cardiac action potentials. In 5th Workshop on Mathematical Foundations of Computational Anatomy, 2015.

http://www.vision.jhu.edu