

Retrieval

Search

I'm Feeling Lucky

offered in: [Français](#)

Plan

1. Introduction
2. Background
3. Image Querying
4. Language Querying
5. Language Subset Querying
6. Conclusion

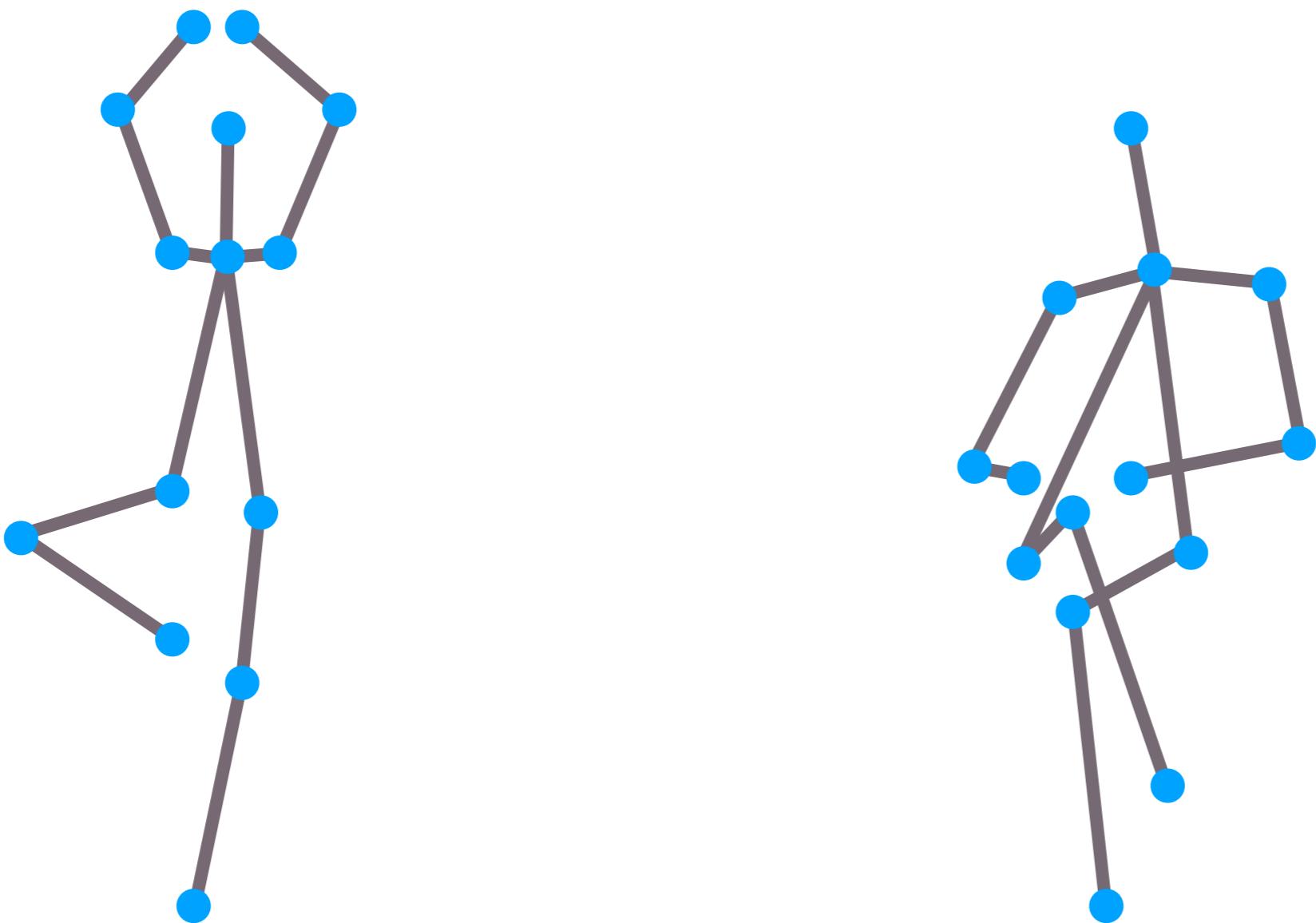


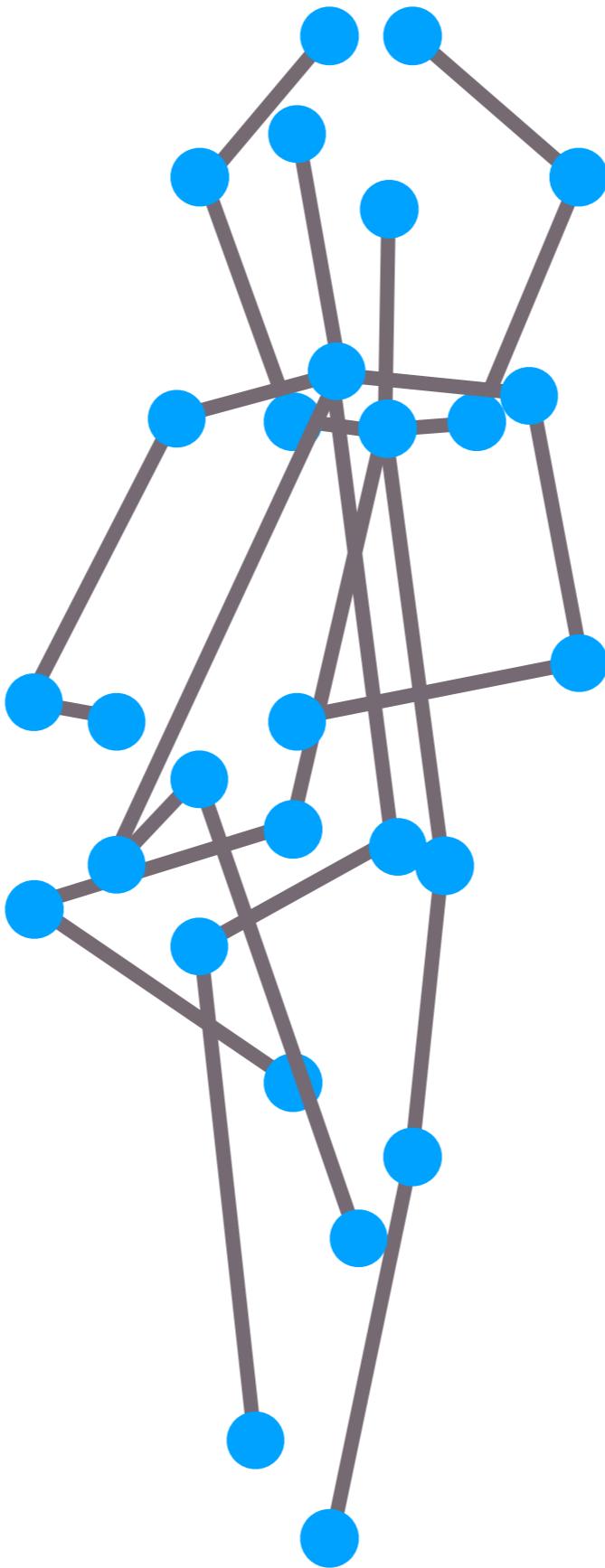


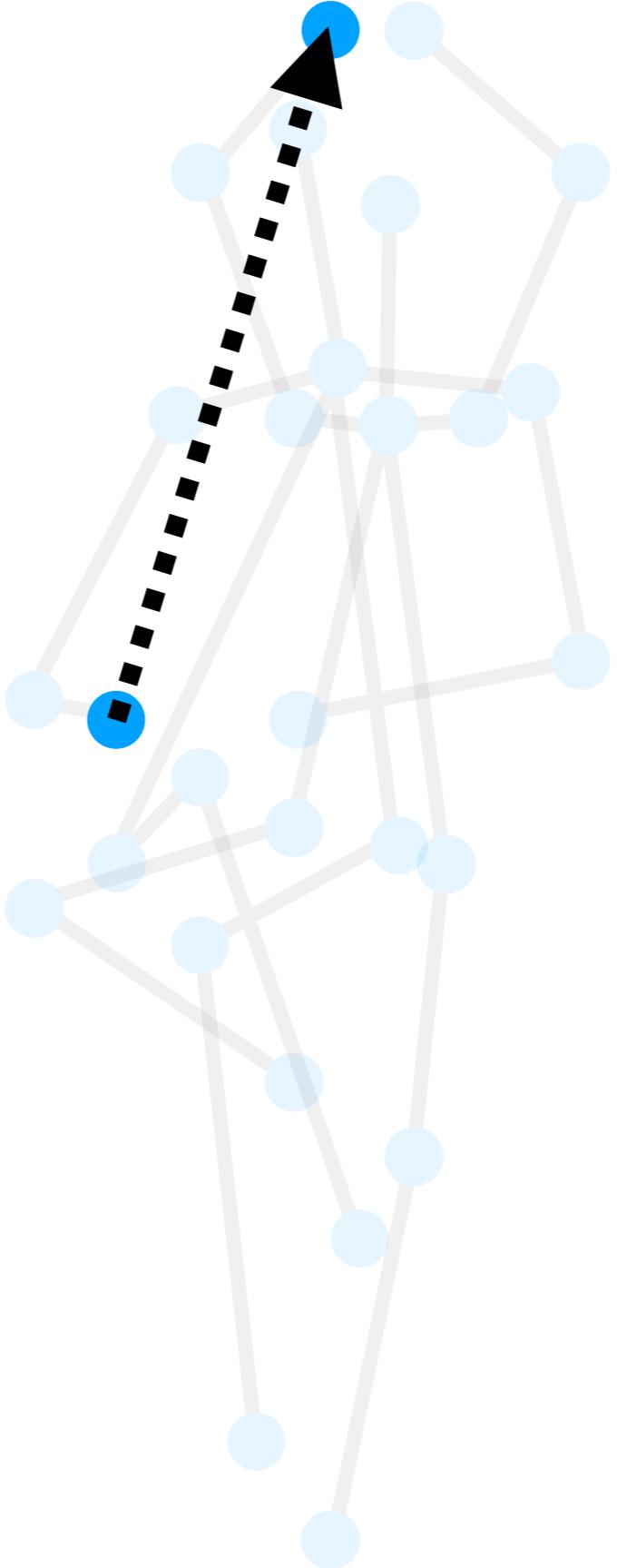












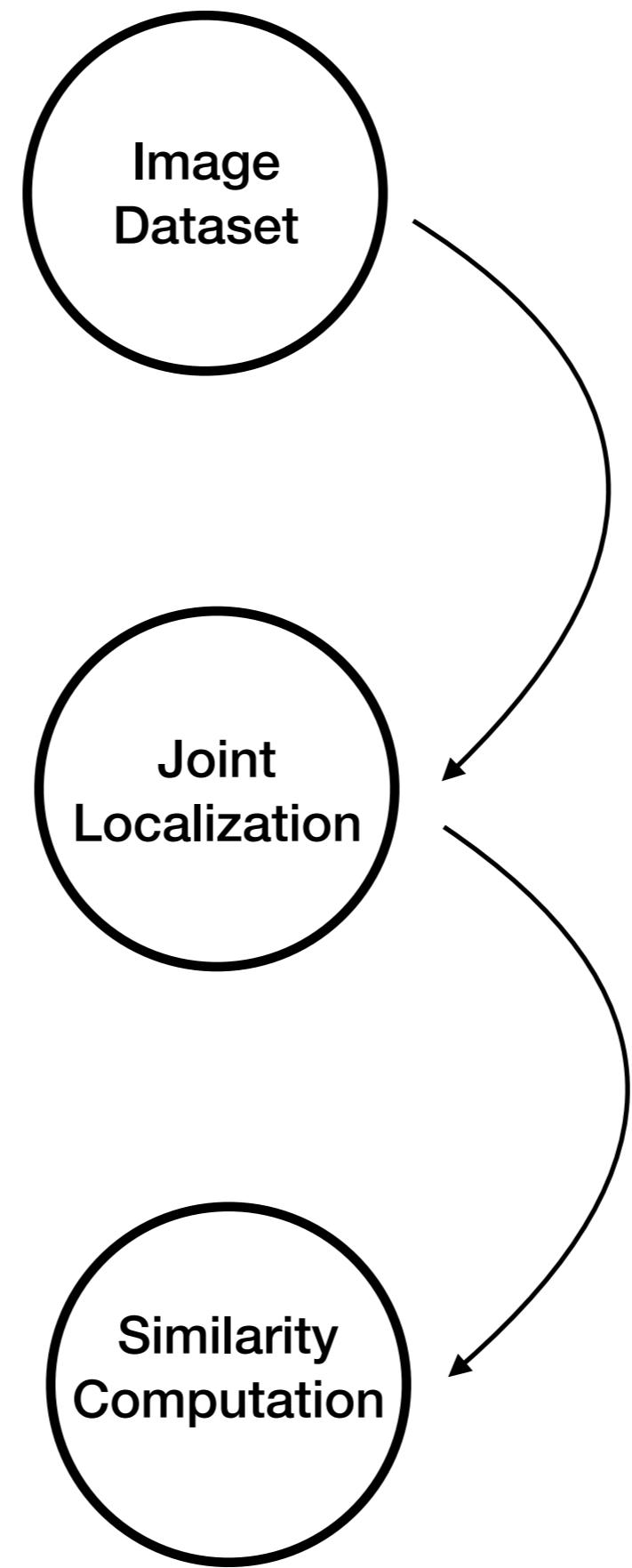


0.1

-

0.13

0.9

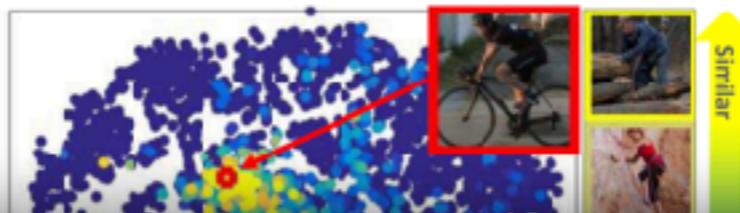


Thin-Slicing for Pose: Learning to Understand Pose without Explicit Pose Estimation

Suha Kwak Minsu Cho Ivan Laptev
 WILLOW Project Team – Inria / École Normale Supérieure, Paris, France

Abstract

We address the problem of learning a pose-aware, compact embedding that projects images with similar human poses to be placed close-by in the embedding space. The embedding is learned from a large dataset of images with various human poses, a real image sequence, and a robust set of images with different human poses by back-projecting them onto a 3D camera rig with a variety of camera settings. The proposed learning framework is able to learn explicit pose embeddings for images with varying body shapes and backgrounds, and it can be applied to a wide range of problems.



Pose Embeddings: A Deep Architecture for Learning to Match Human Poses

1. Intr

Greg Mori^{*1}, Caroline Pantofaru², Nisarg Kothari², Thomas Leung², George Toderici², Alexander Toshev², and Weilong Yang²

¹Simon Fraser University

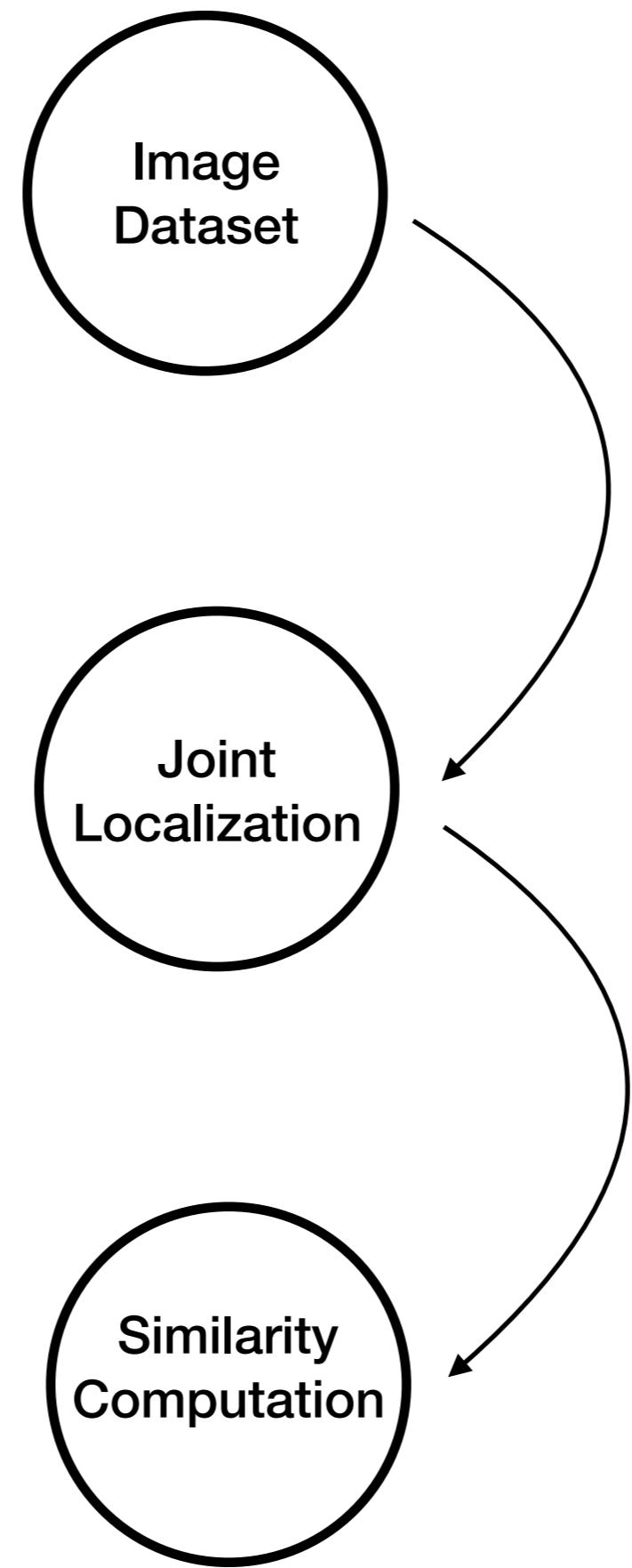
²Google Inc.

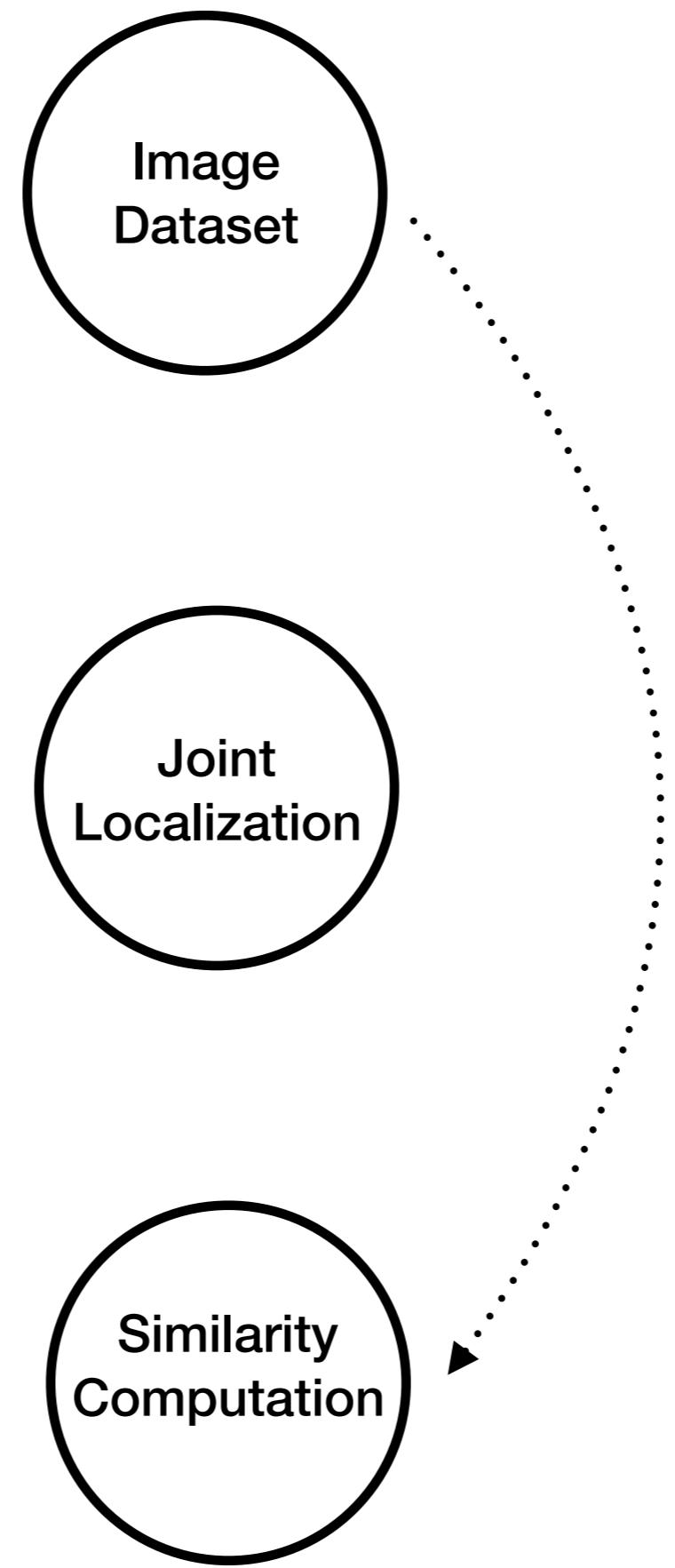
Abstract

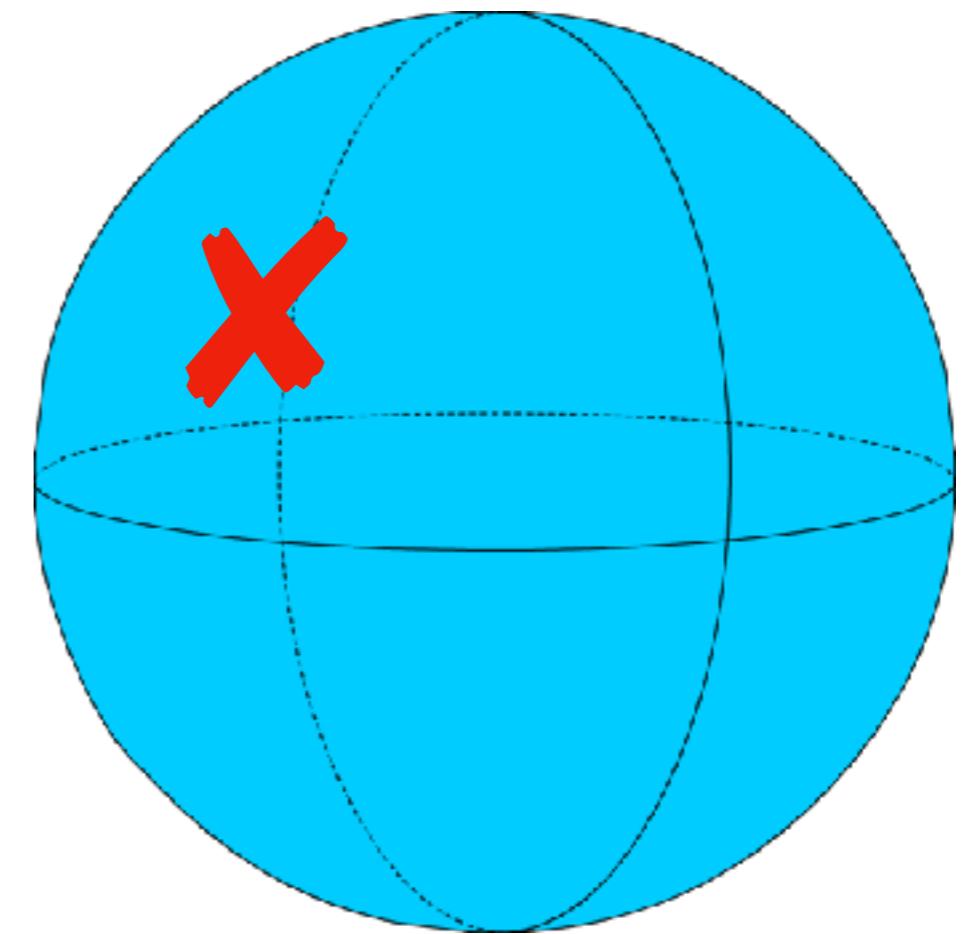
We present a method for learning an embedding that places images of humans in similar poses nearby. This embedding can be used as a direct method of comparing im-

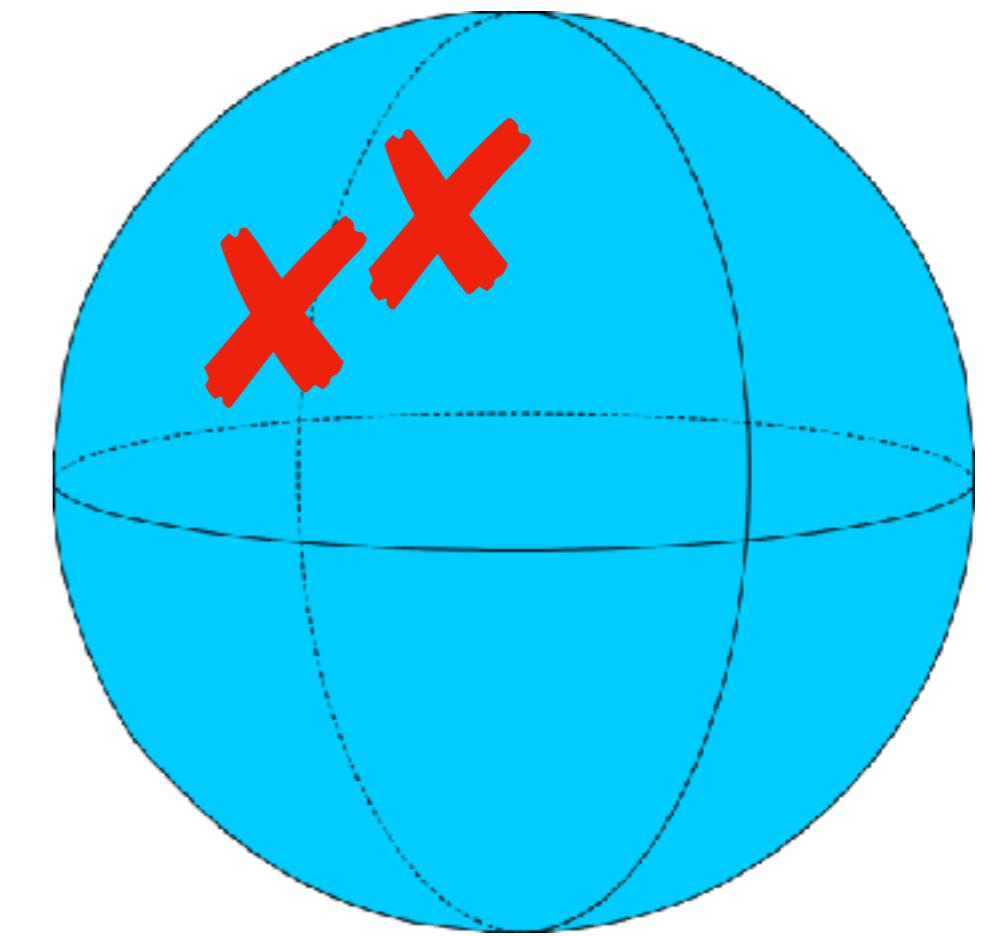


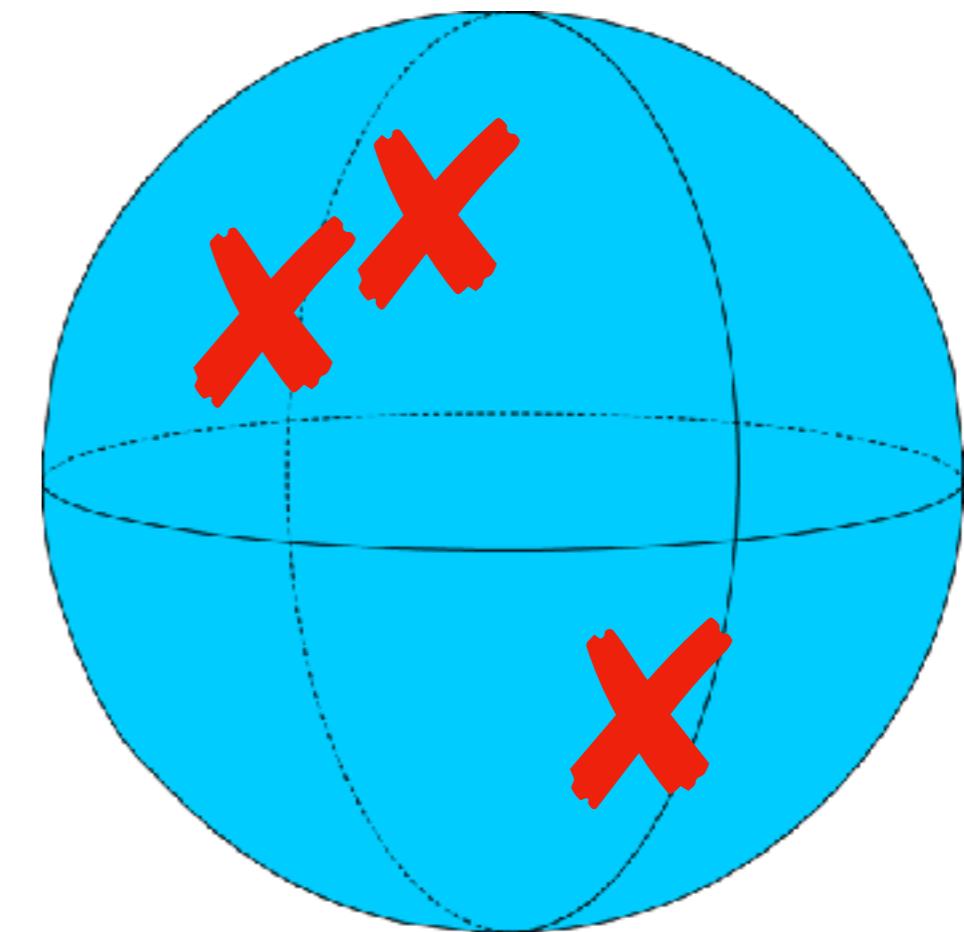
=?

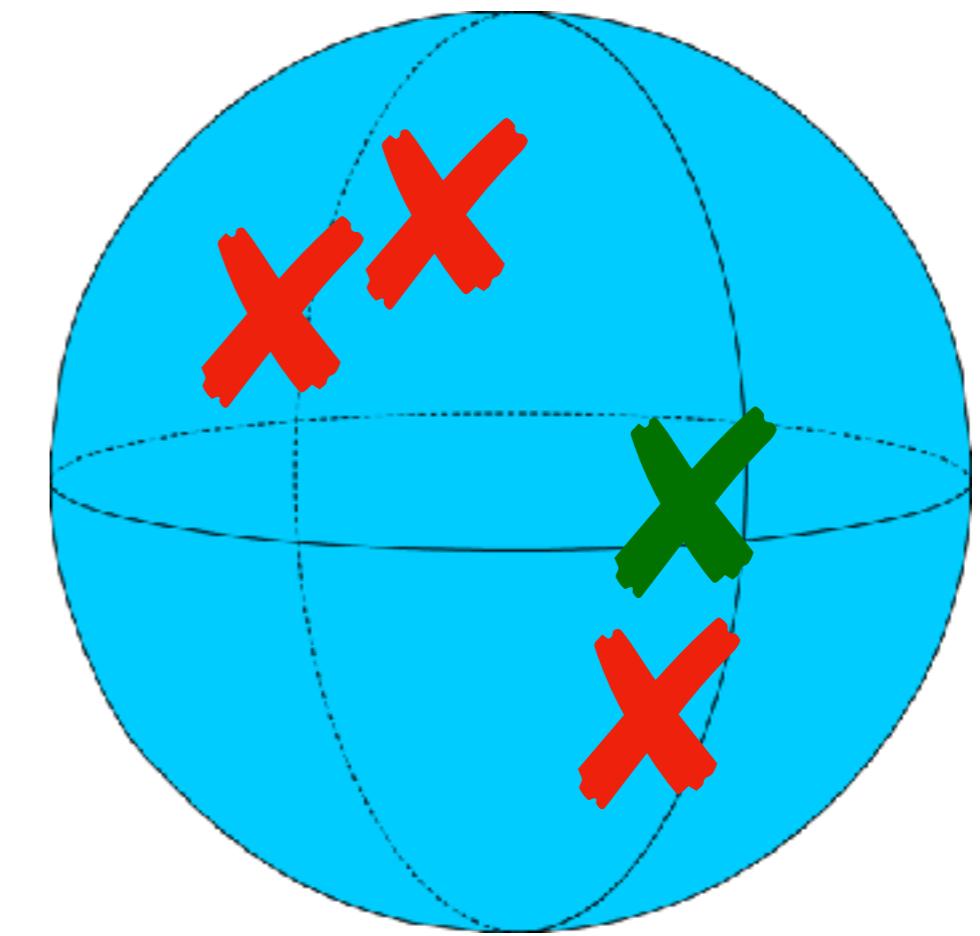


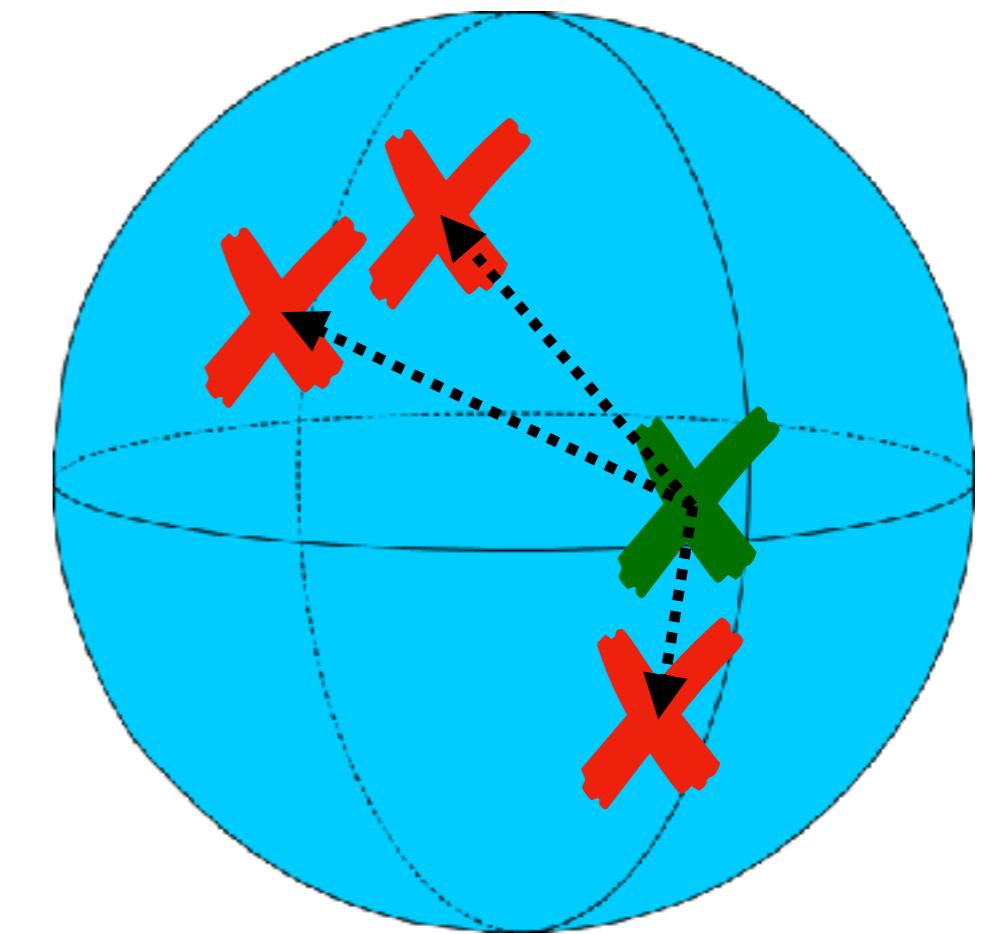
















Main Contributions

Pose-Aware embedding networks

Image-Language embedding network

Two solutions for language subsets

A dataset of 2D, 3D, and Language
Primitive descriptors

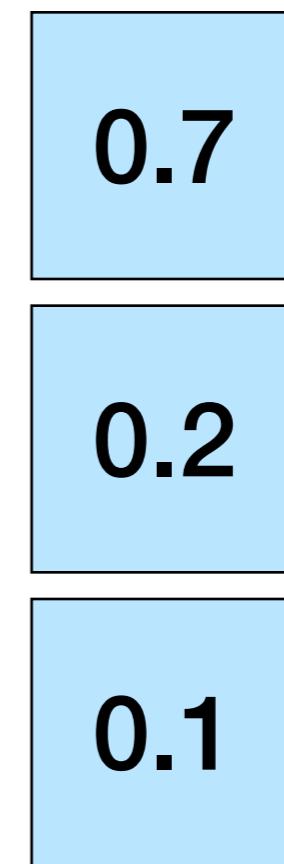
Background



Baseball

Tennis

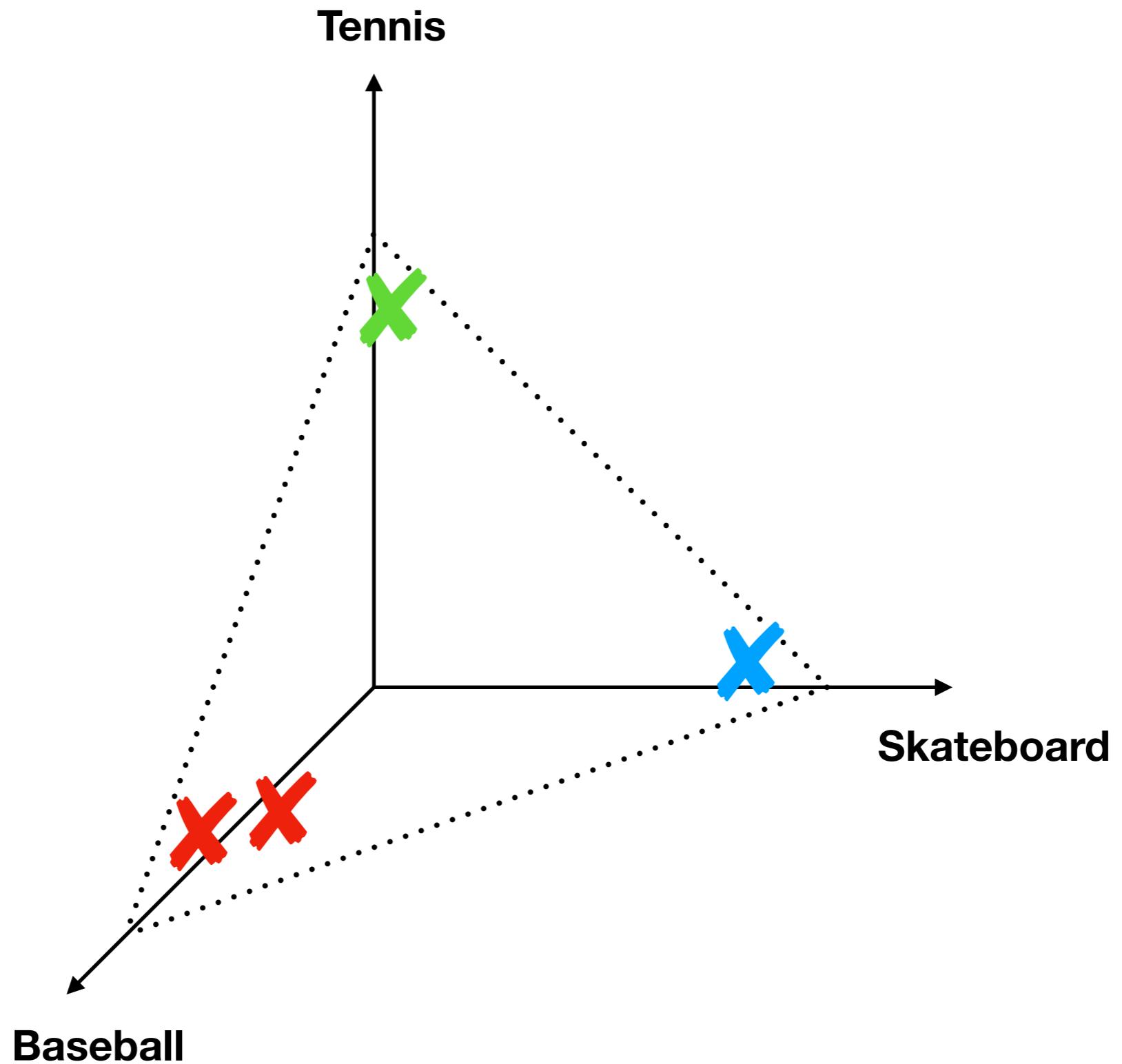
Skateboard



Baseball

Tennis

Skateboard





x

y

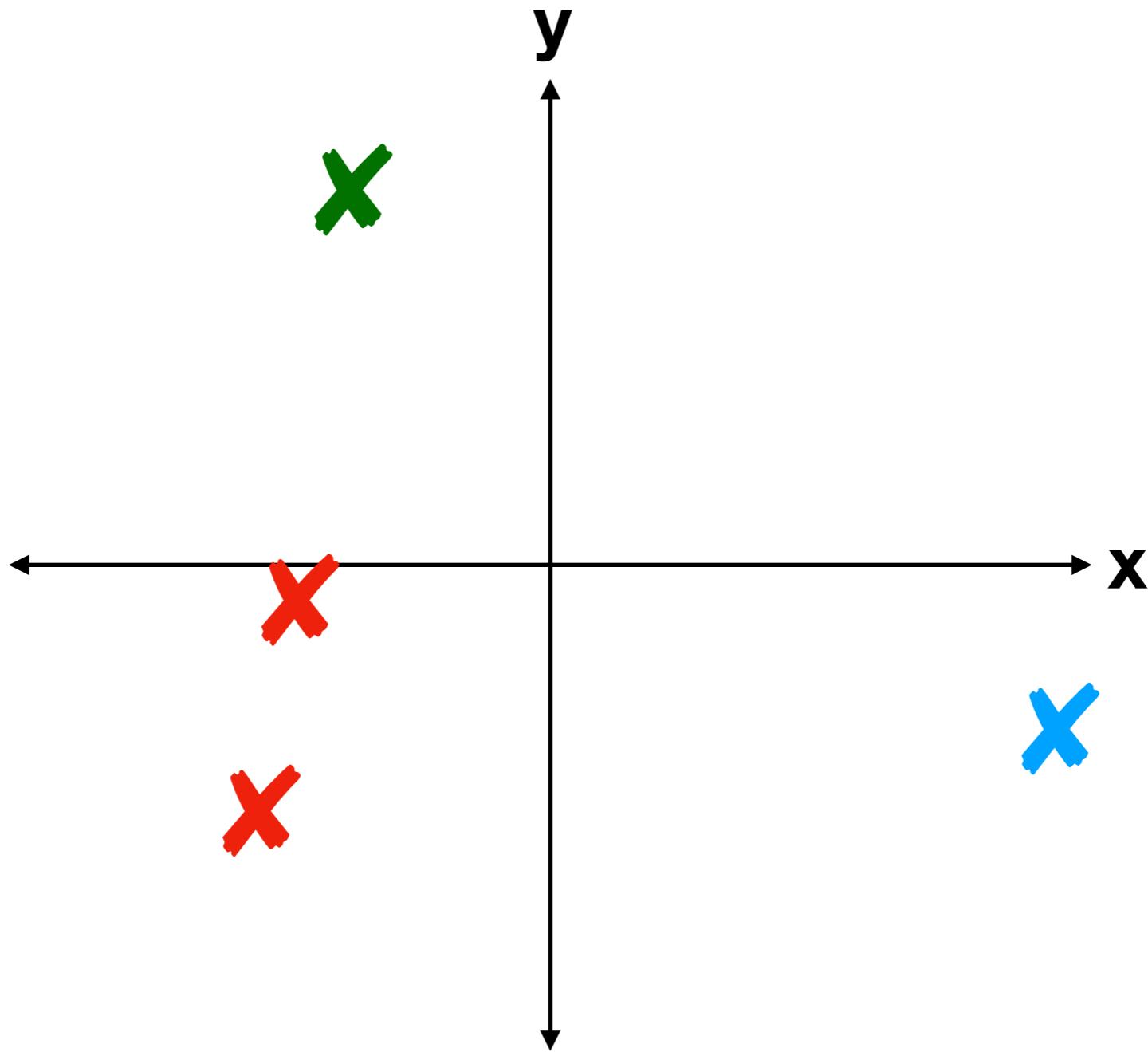


0.3

0.9

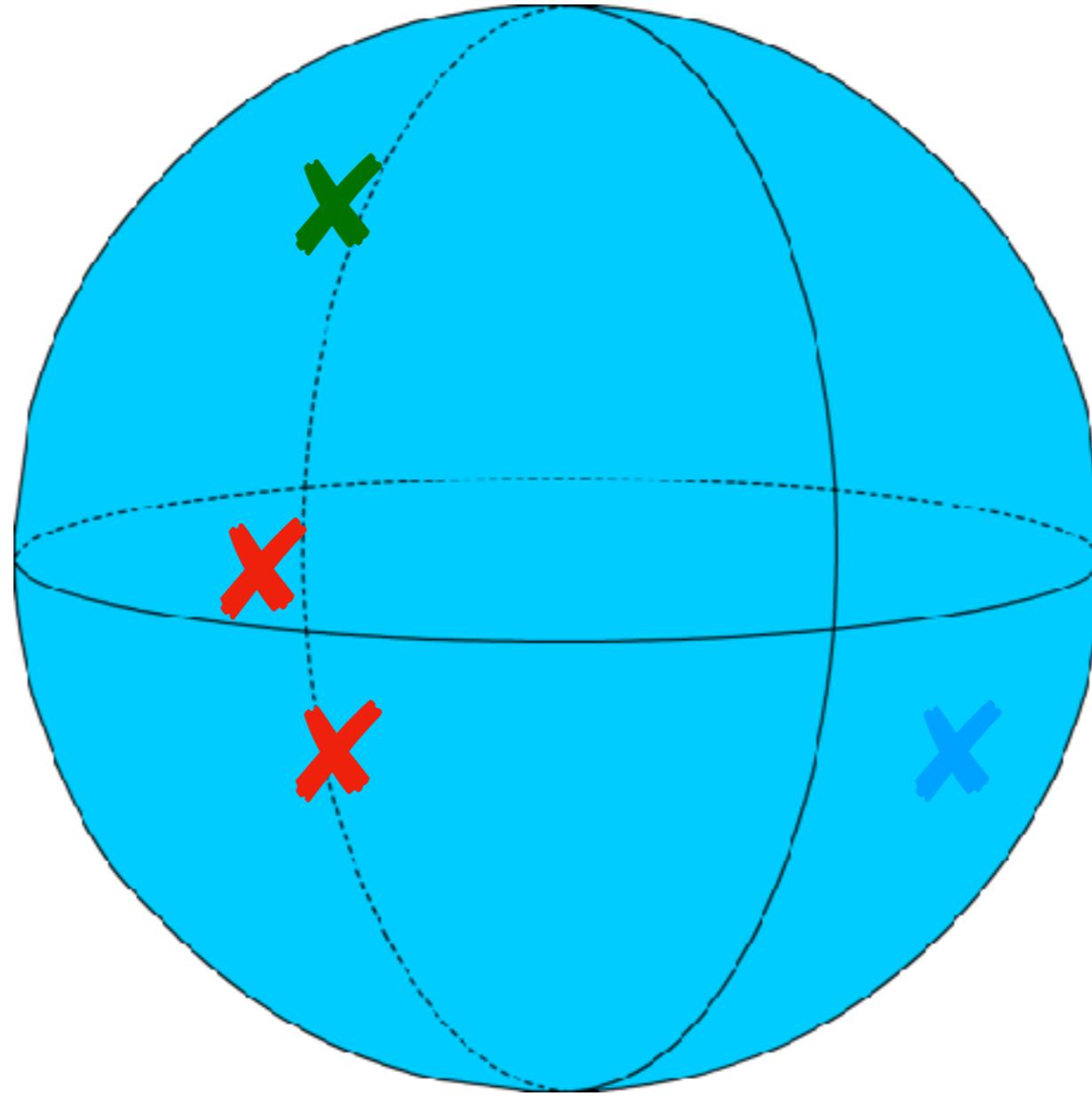
x

y



$$\|f(a) - f(p)\|_2^2 < \|f(a) - f(n)\|_2^2$$

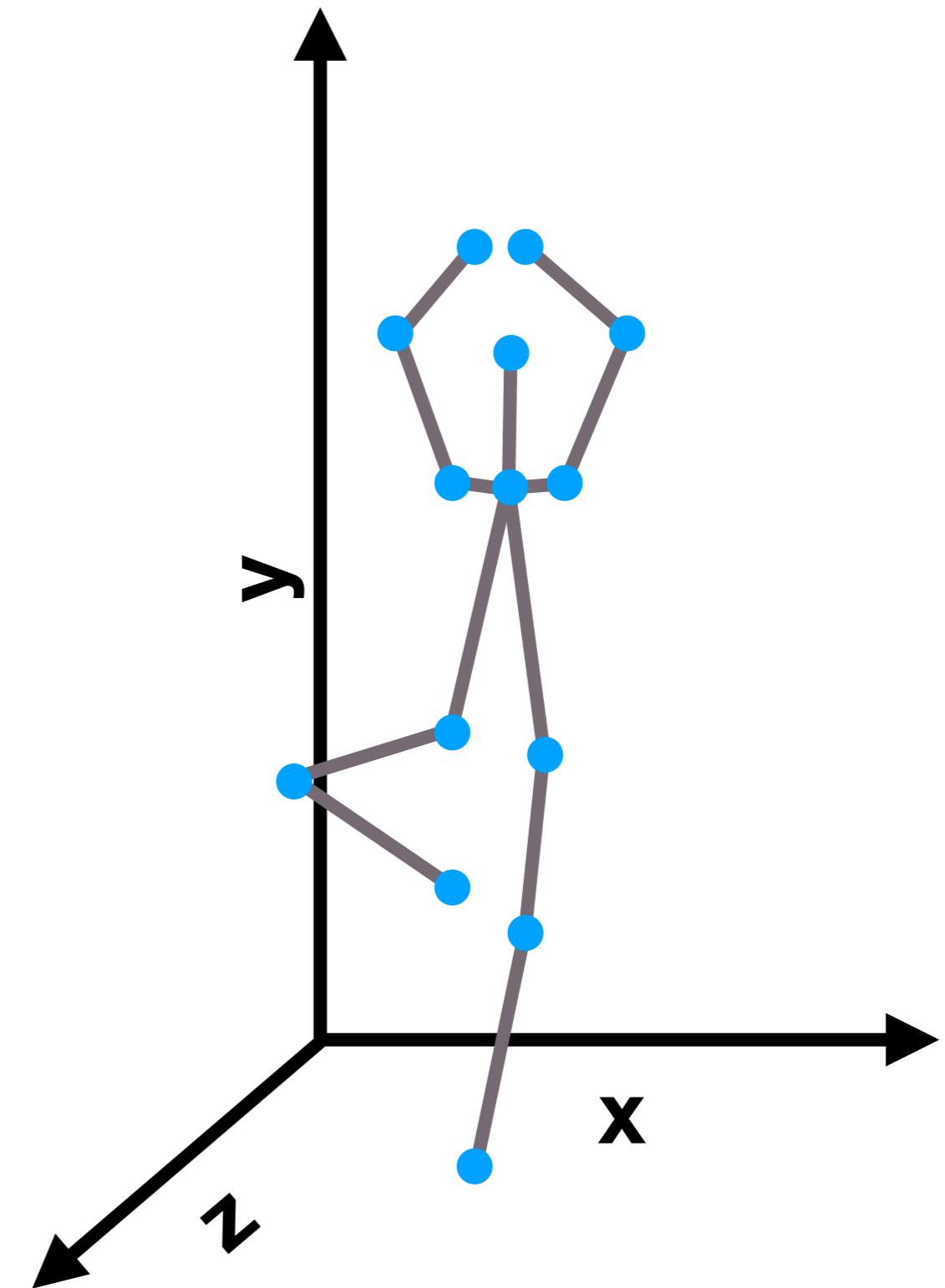
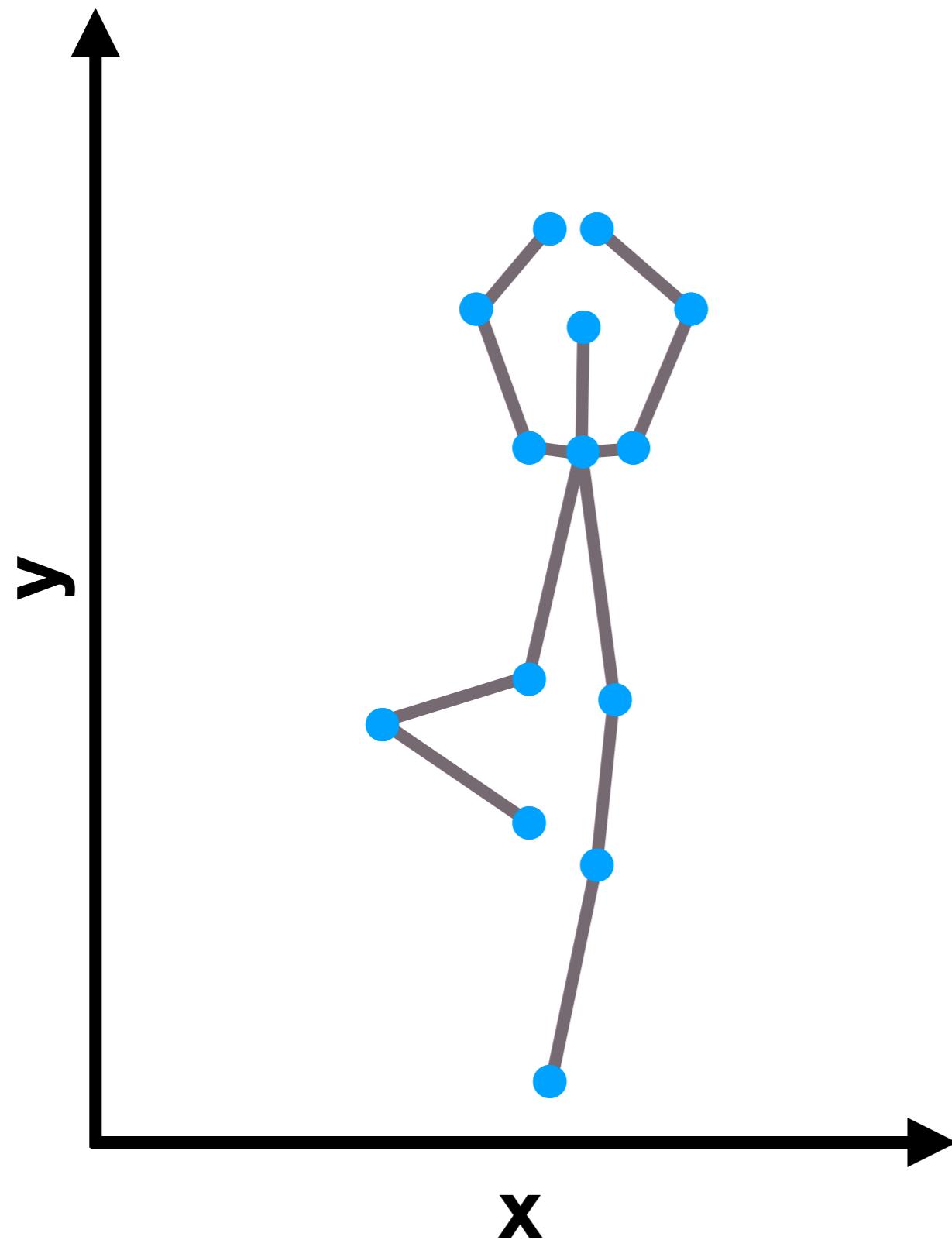
$\forall a, p, n \in D$ s.t. $\delta(a_{pose}, p_{pose}) < \delta(a_{pose}, n_{pose})$

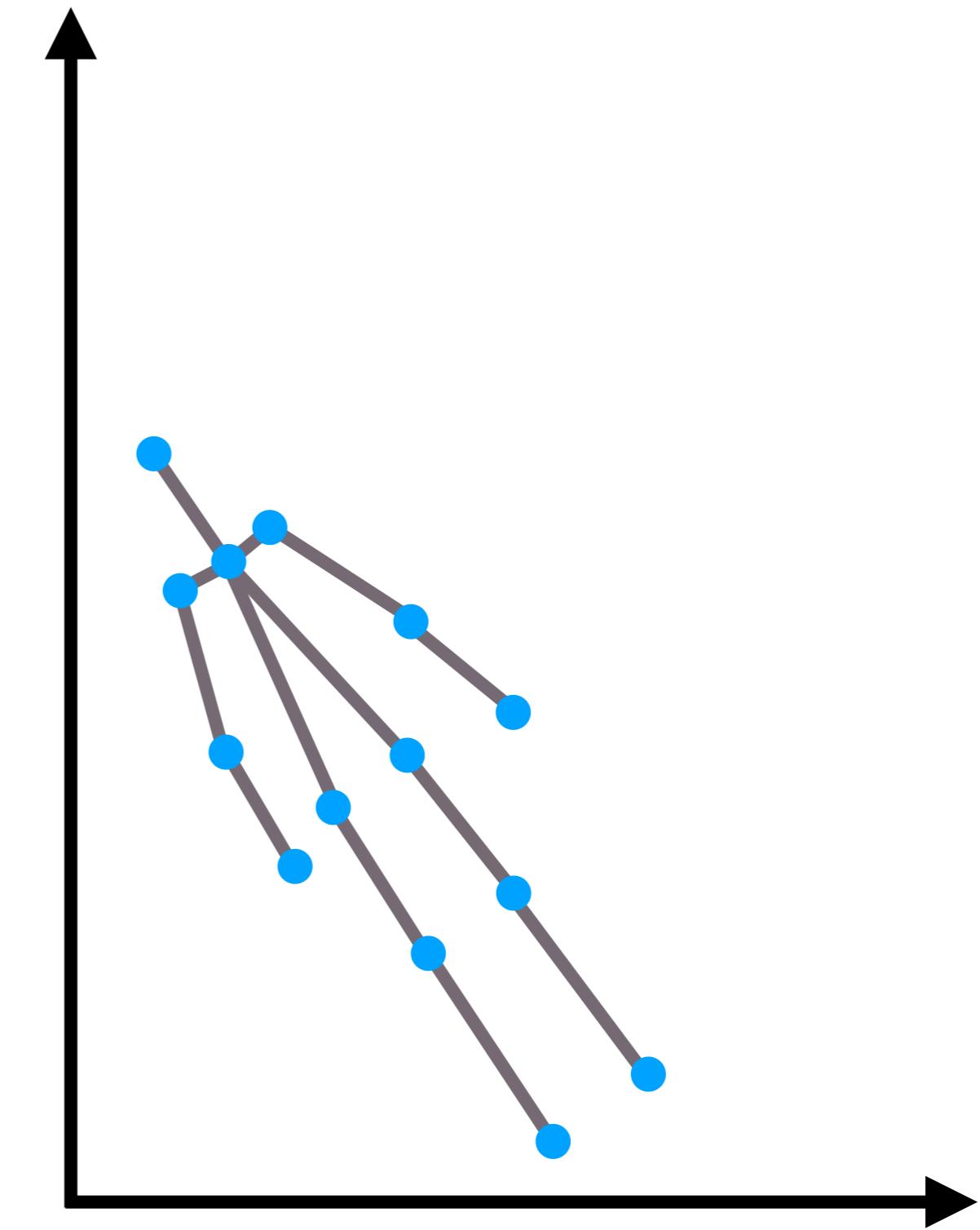
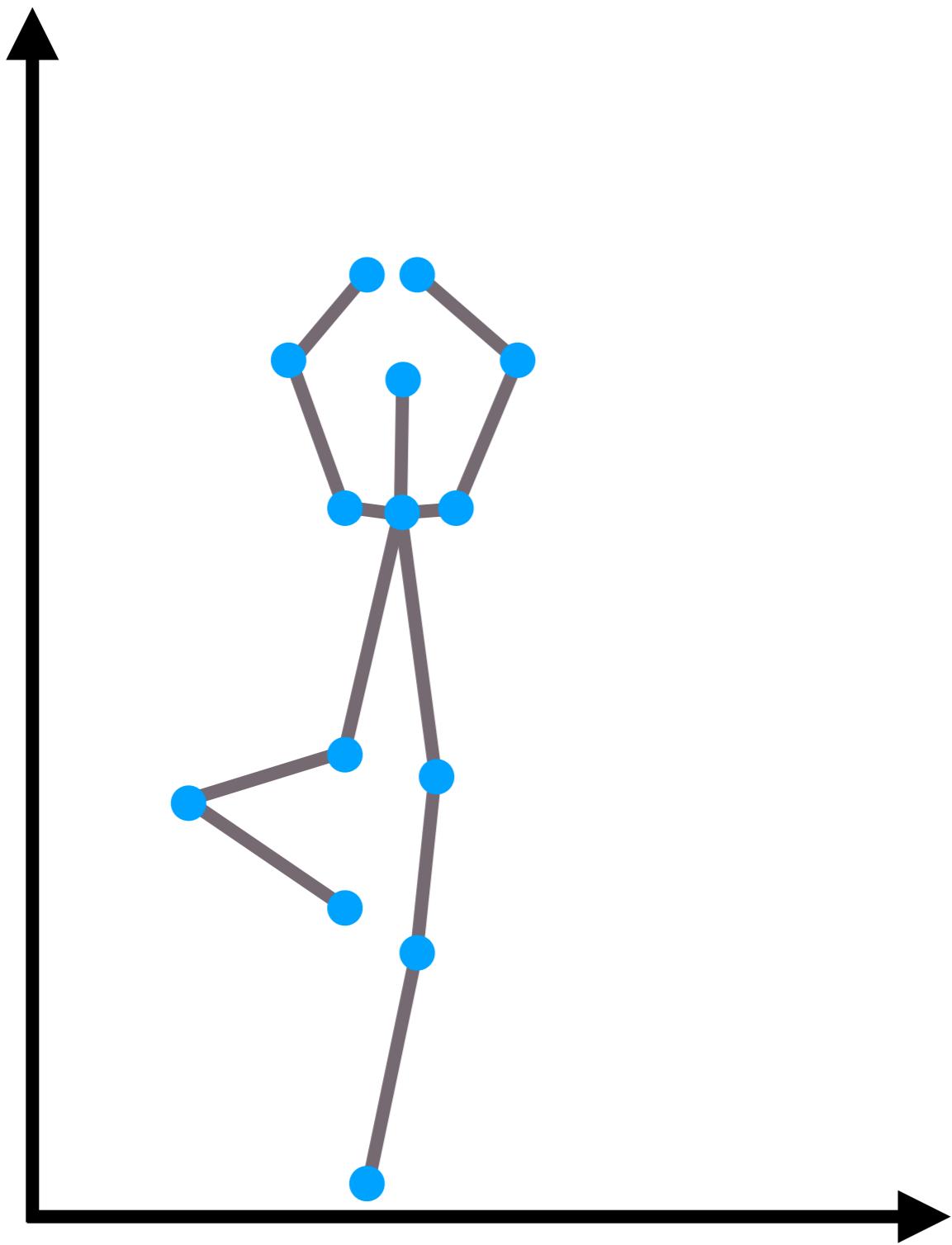


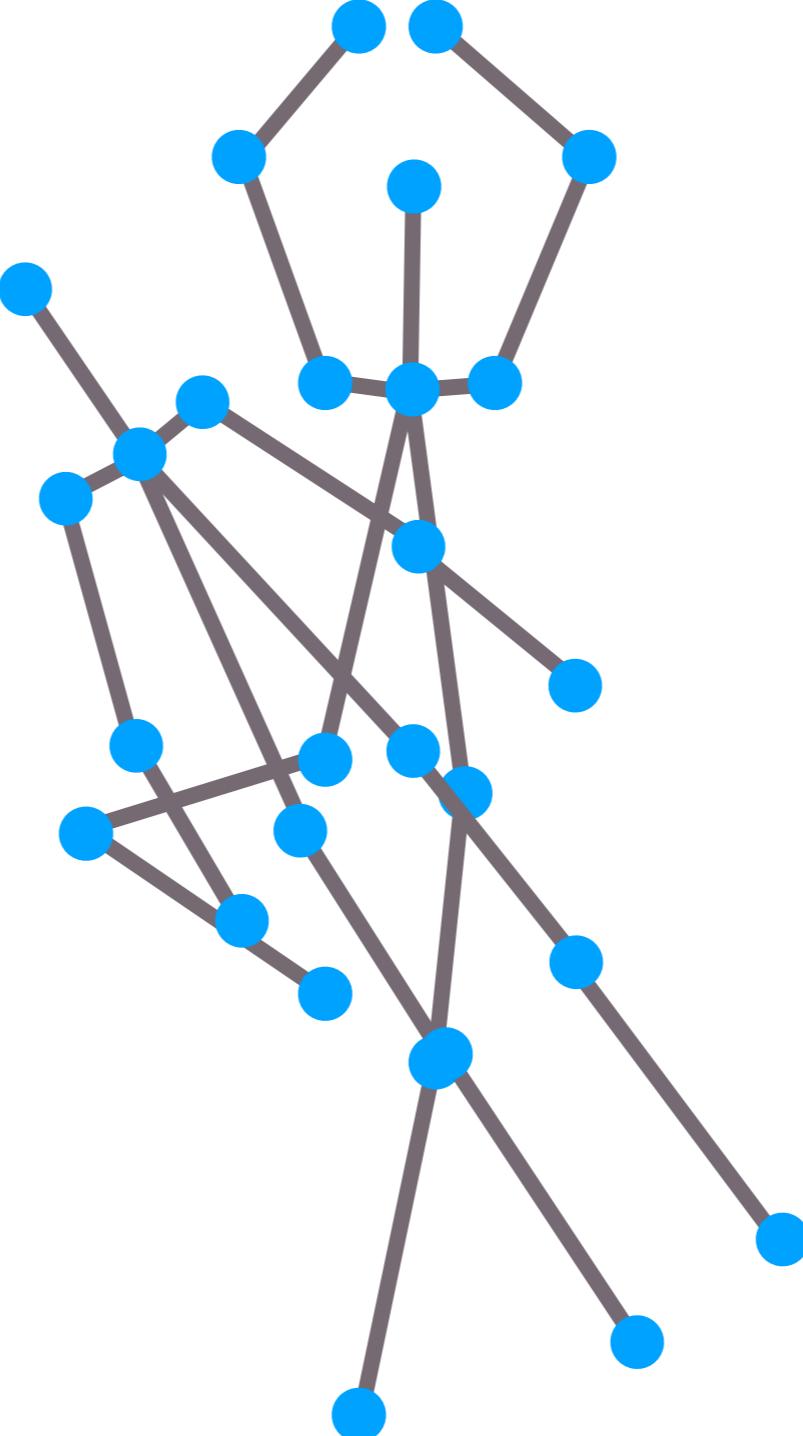
$$\|f(a) - f(p)\|_2^2 < \|f(a) - f(n)\|_2^2$$

$\forall a, p, n \in D$ s.t. $\delta(a_{pose}, p_{pose}) < \delta(a_{pose}, n_{pose})$

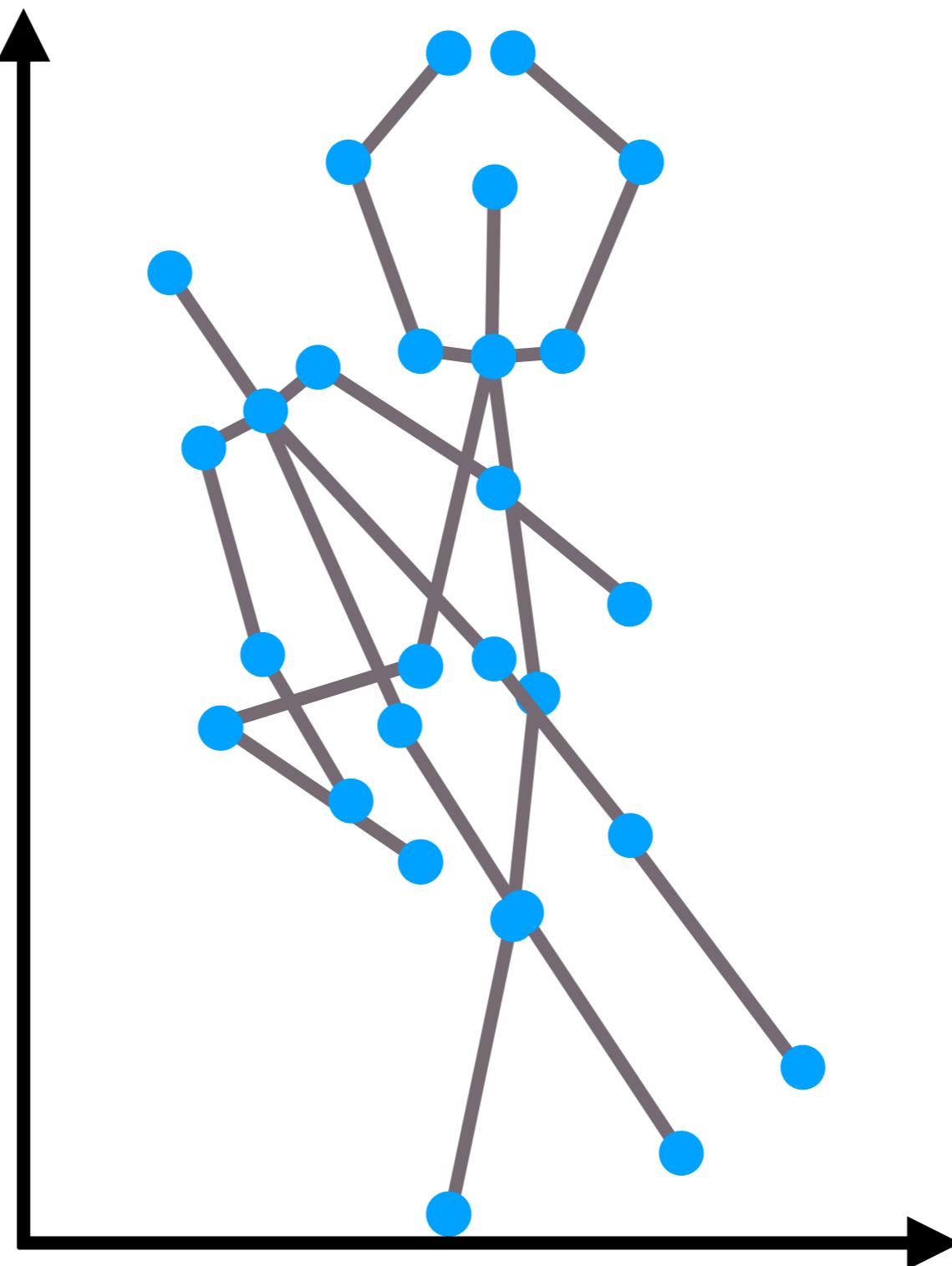
Image Querying



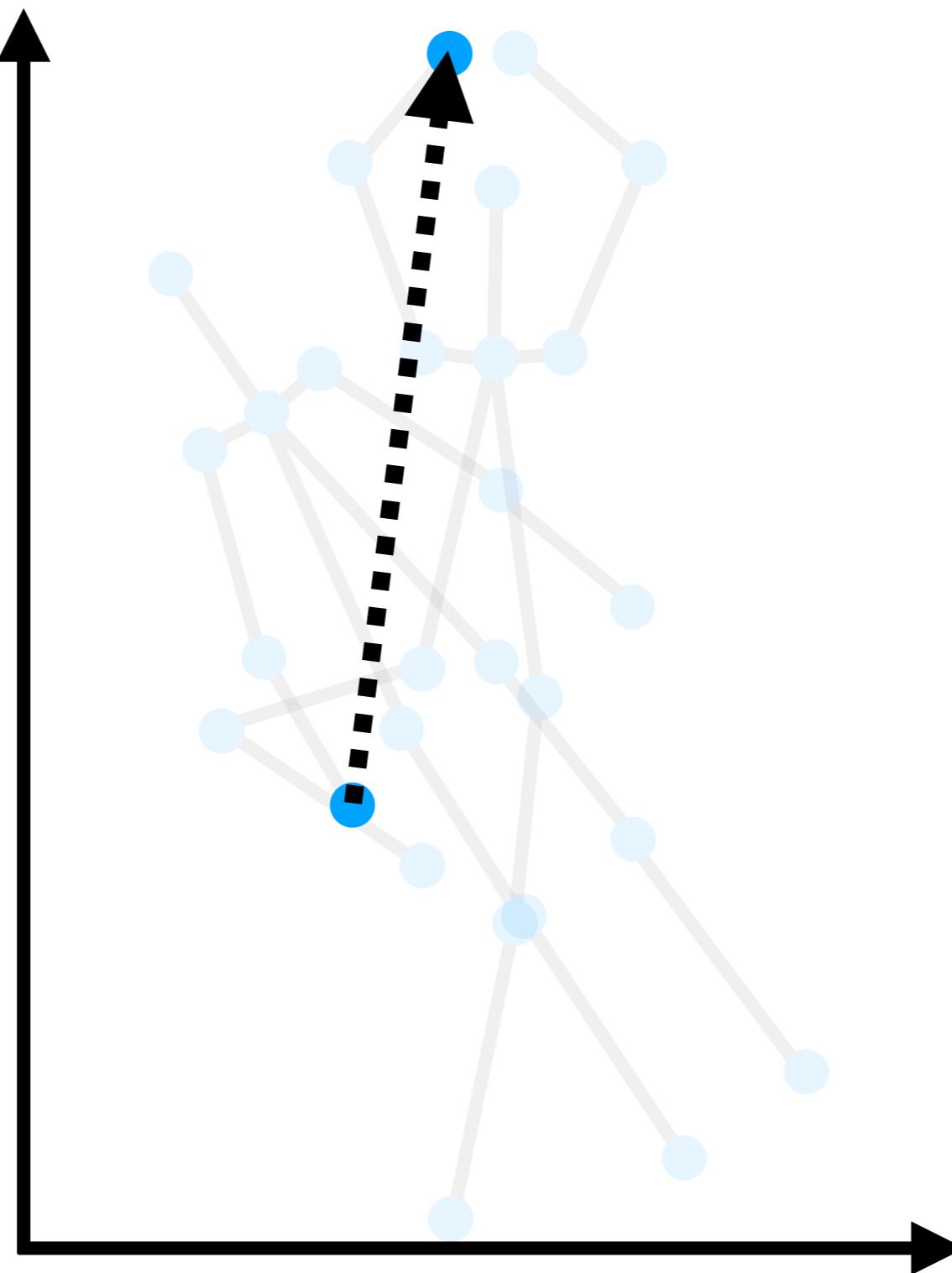




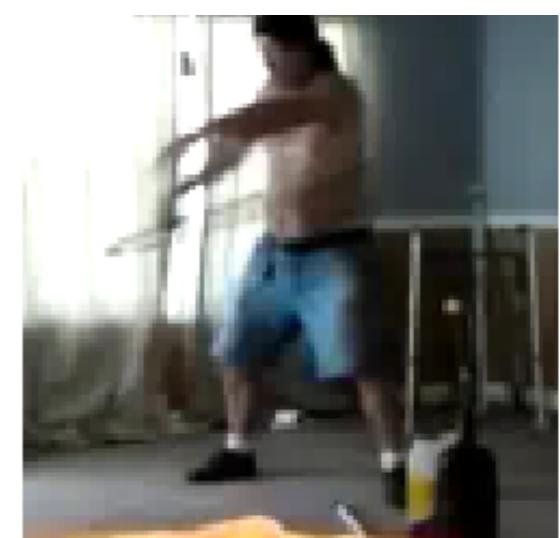
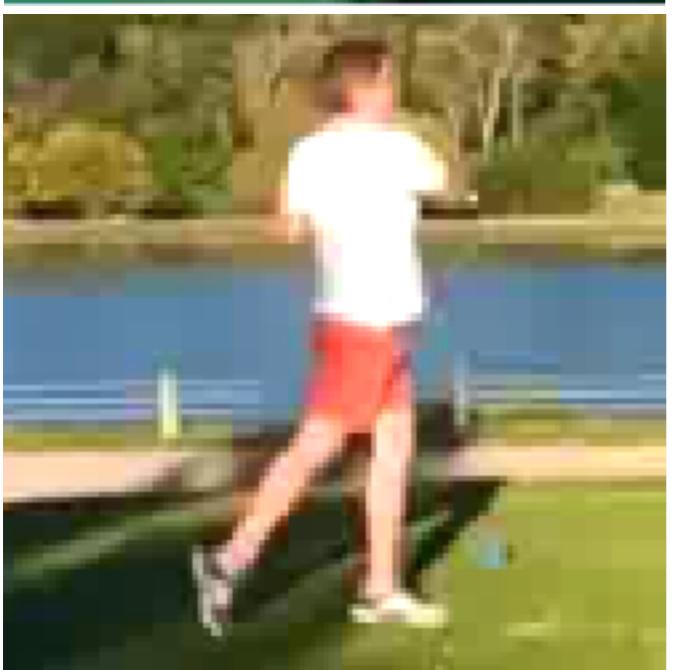
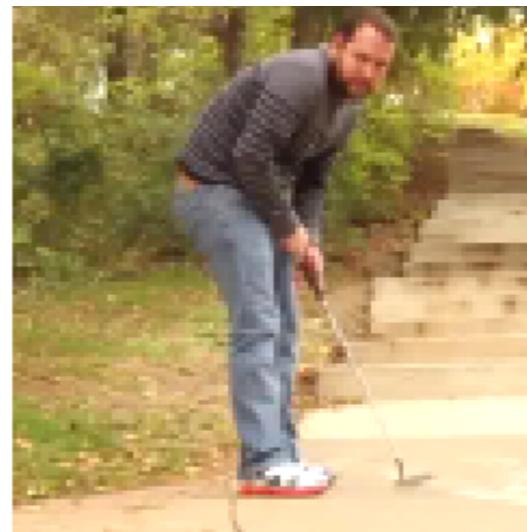
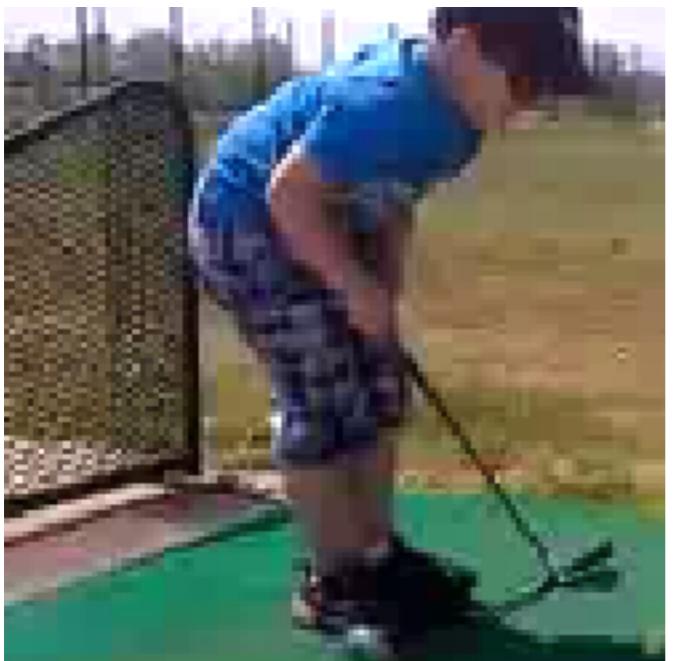
$$\delta_{2D}(J_a, J_b) = \frac{1}{N} \sum_{i=1}^N \sqrt{|J_a^i - J_b^i|^2}$$

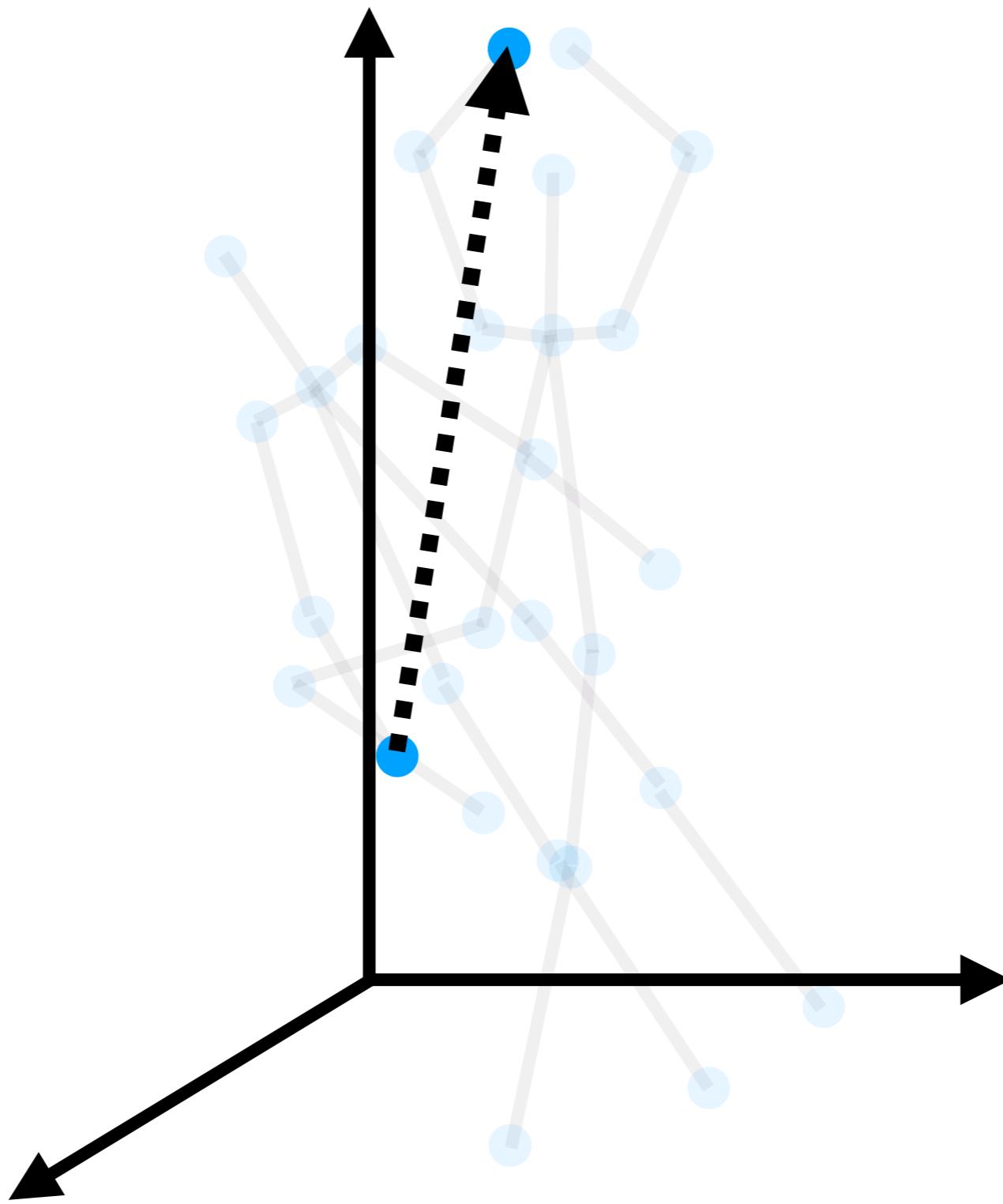


$$\delta_{2D}(J_a, J_b) = \frac{1}{N} \sum_{i=1}^N \sqrt{|J_a^i - J_b^i|^2}$$

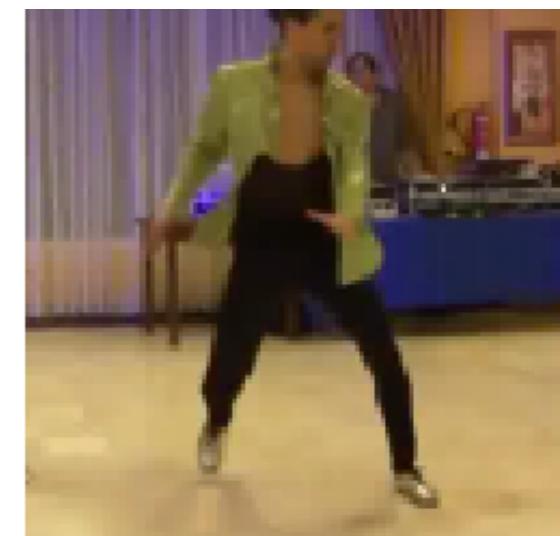
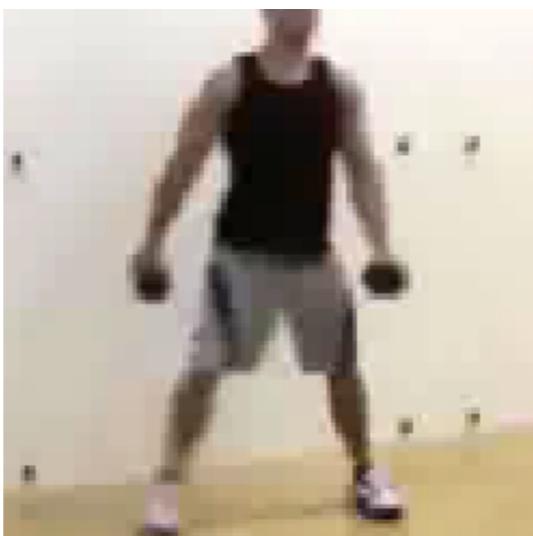
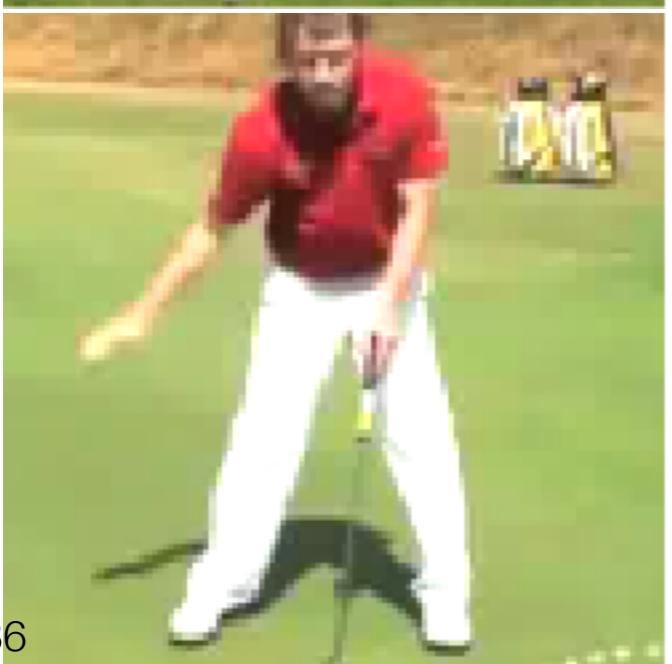
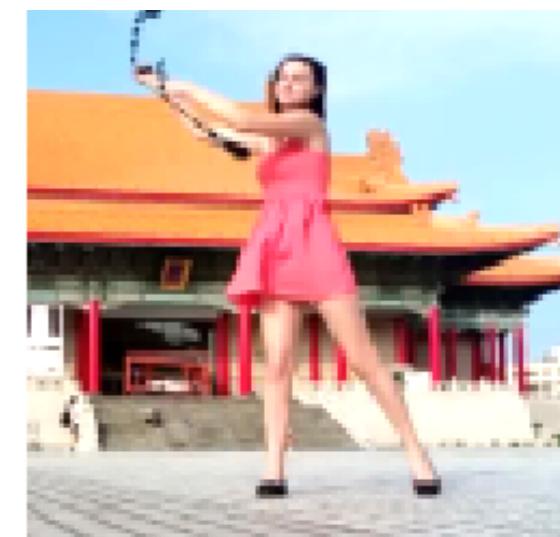
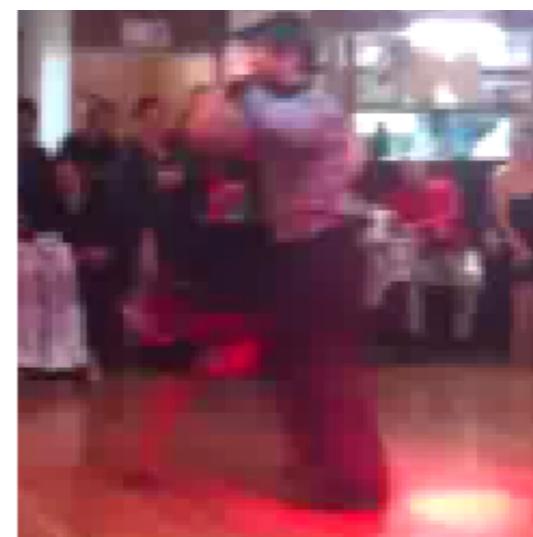
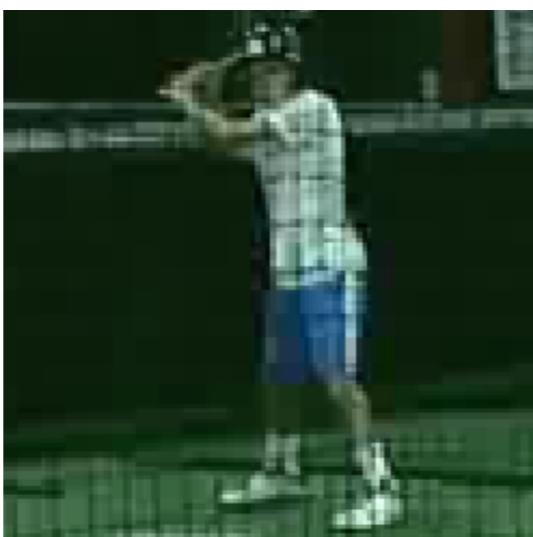
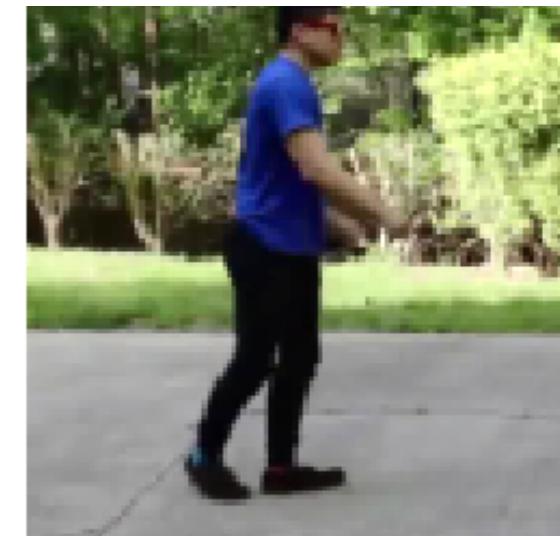


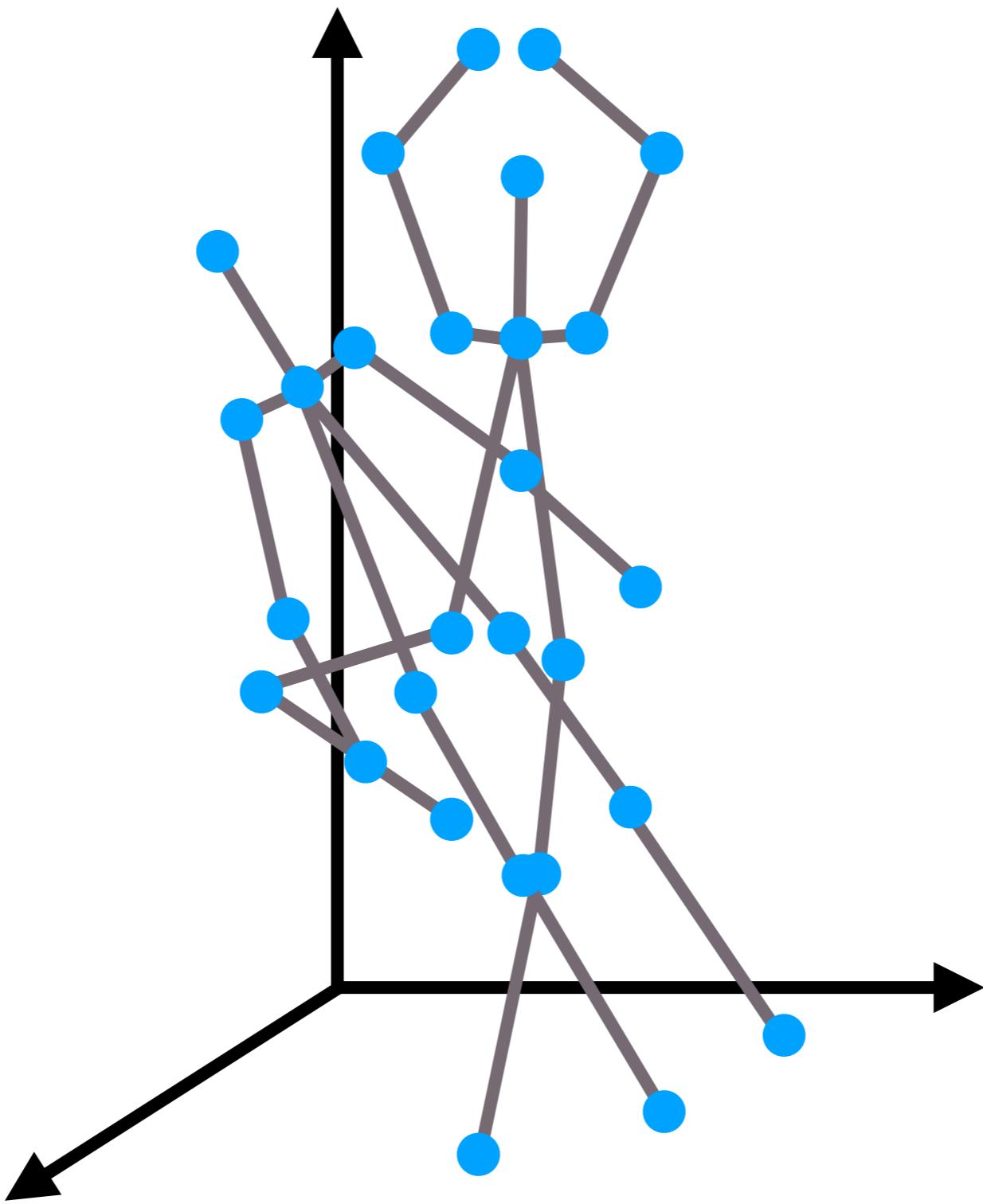
$$\delta_{2D}(J_a, J_b) = \frac{1}{N} \sum_{i=1}^N \sqrt{|J_a^i - J_b^i|^2}$$

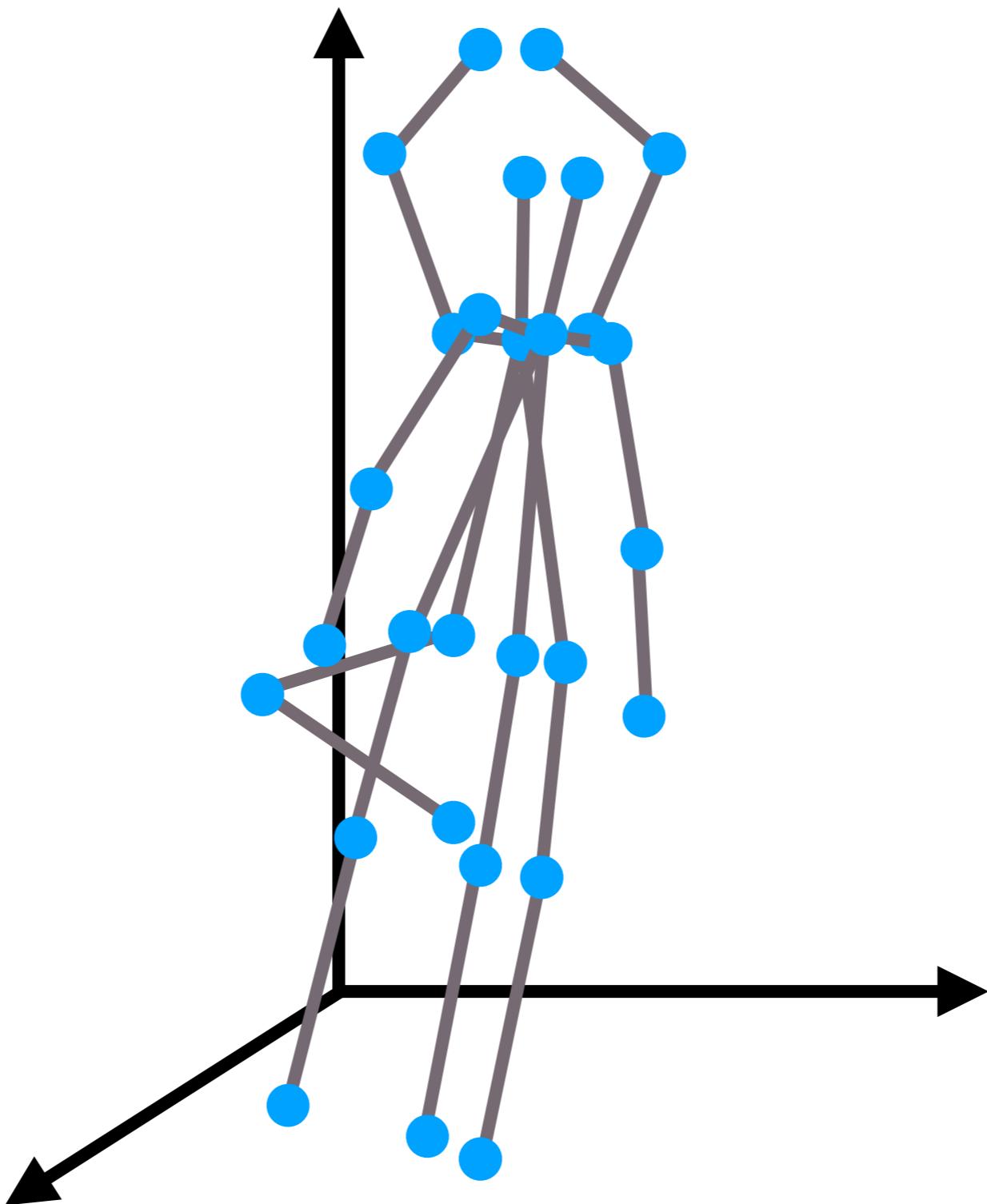




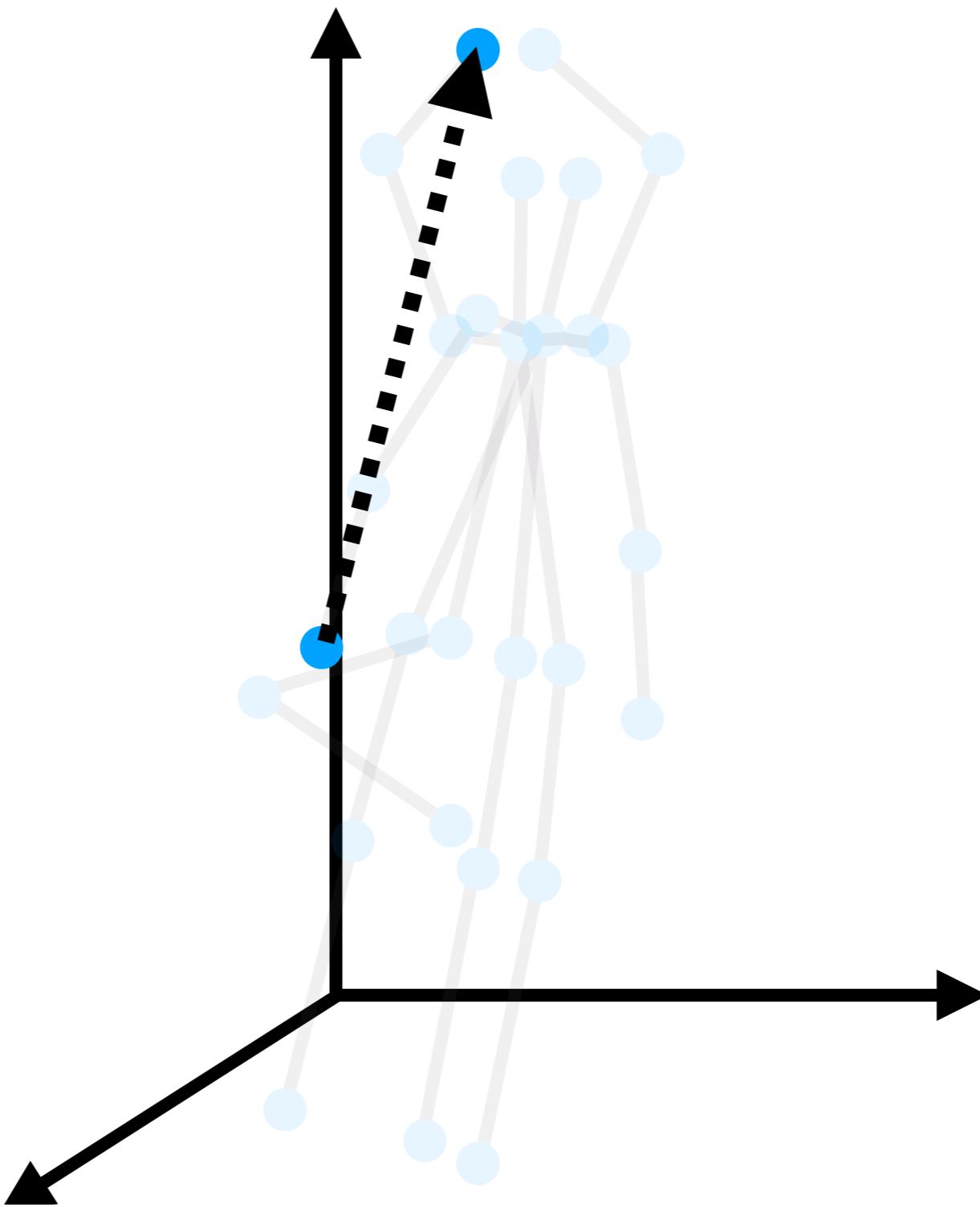
$$\delta_{3D}(J_a, J_b) = \frac{1}{N} \sum_{i=1}^N \sqrt{|J_a^i - J_b^i|^2}$$



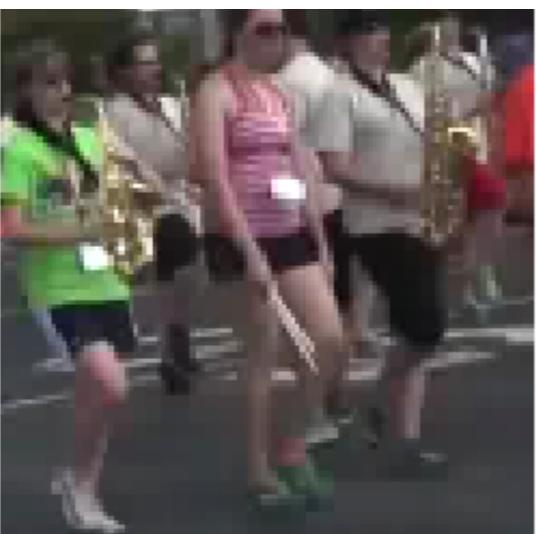
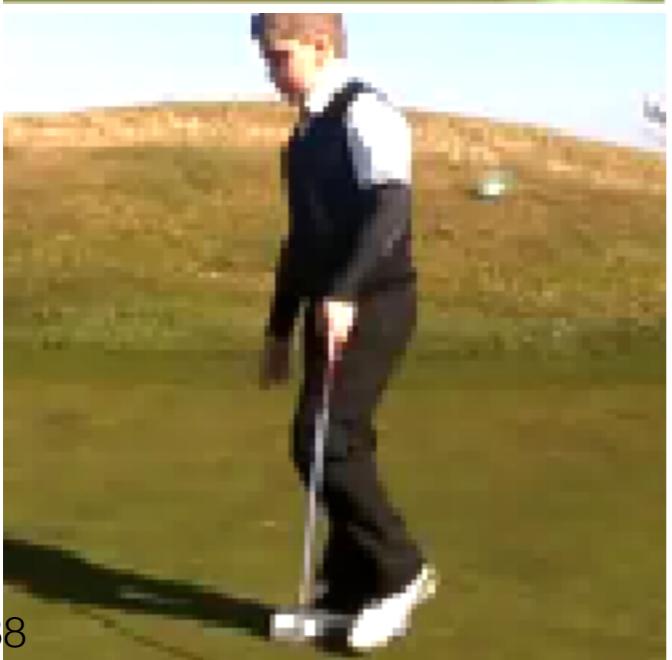
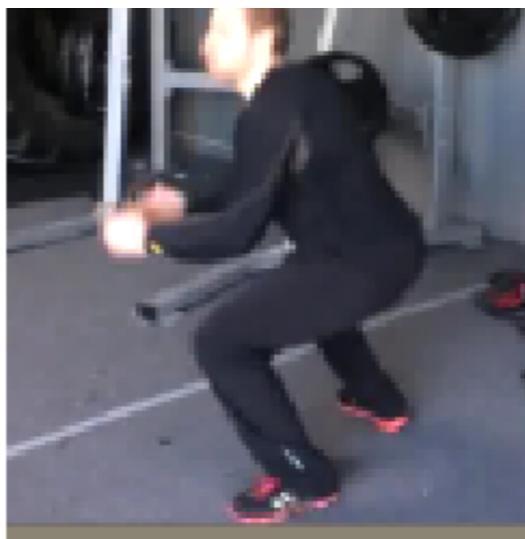
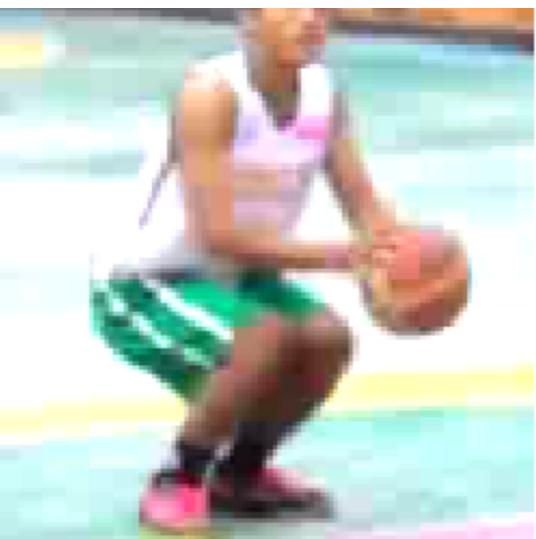
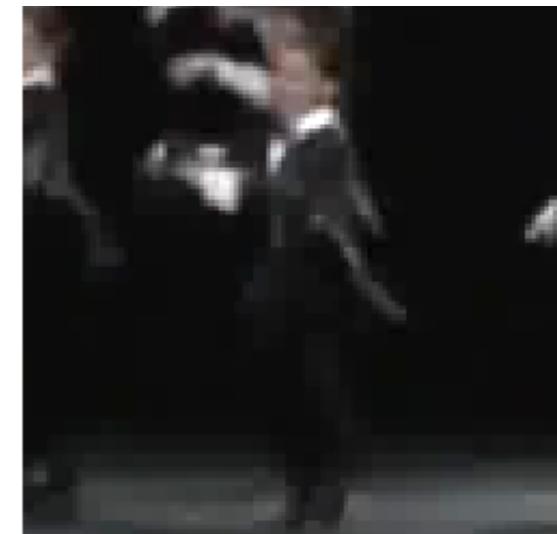
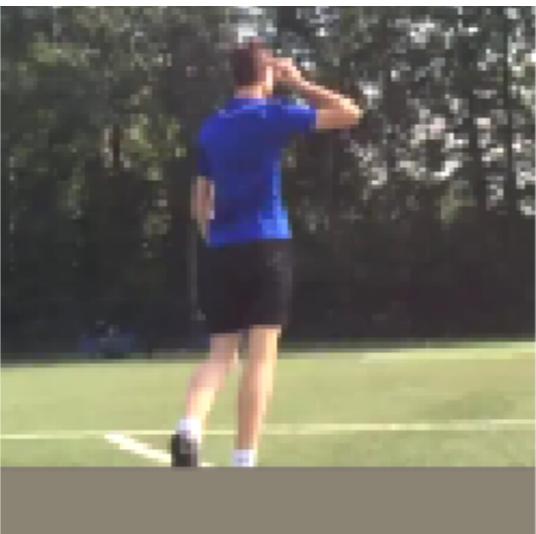
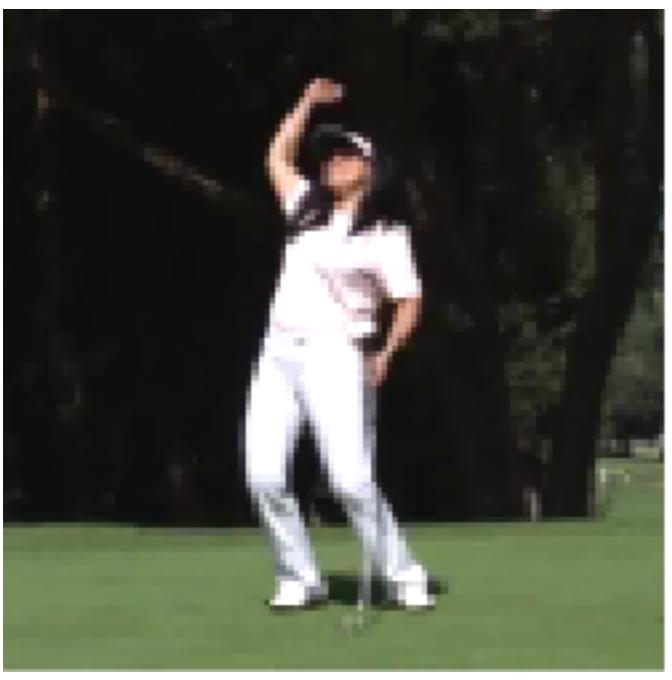


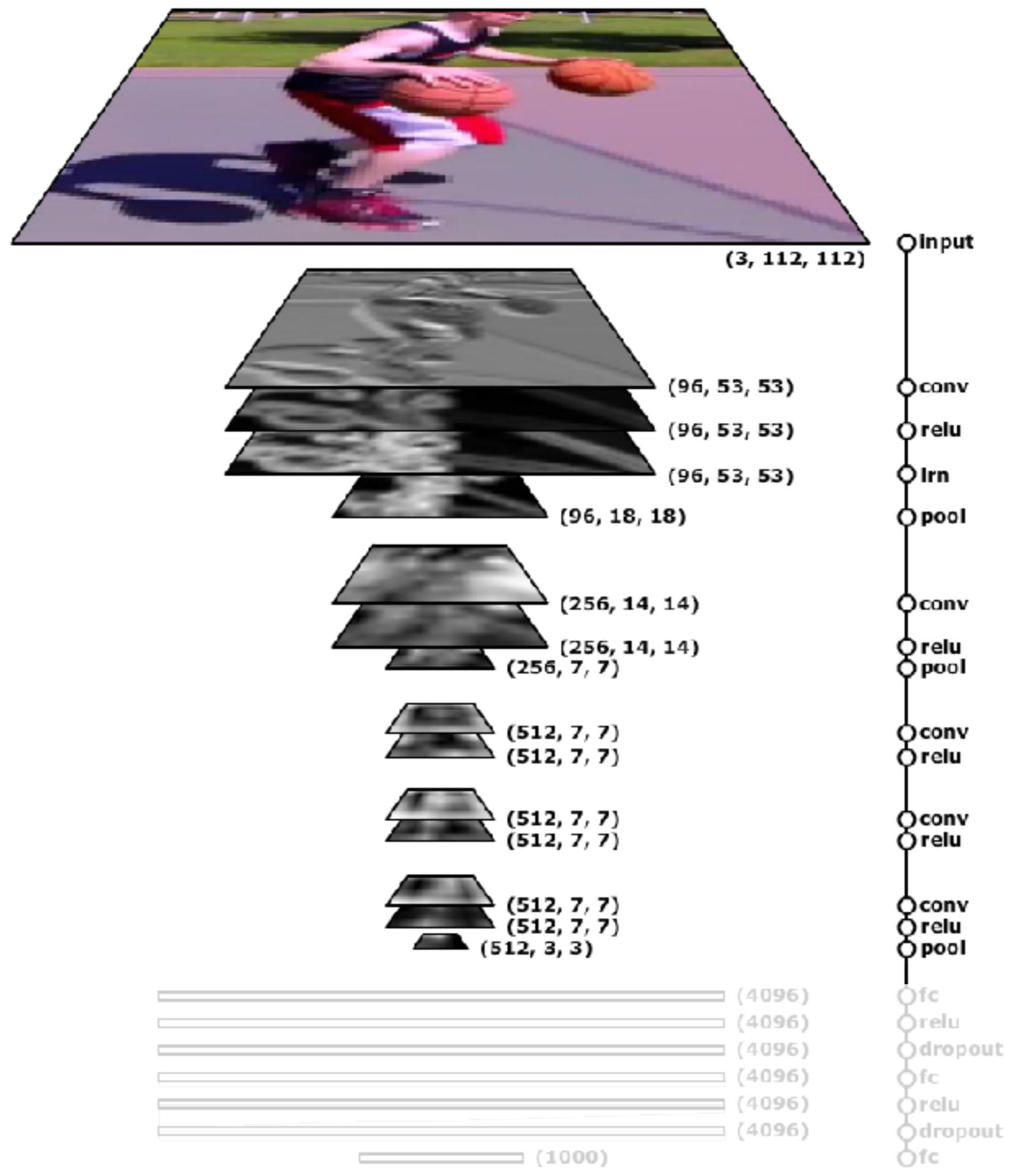


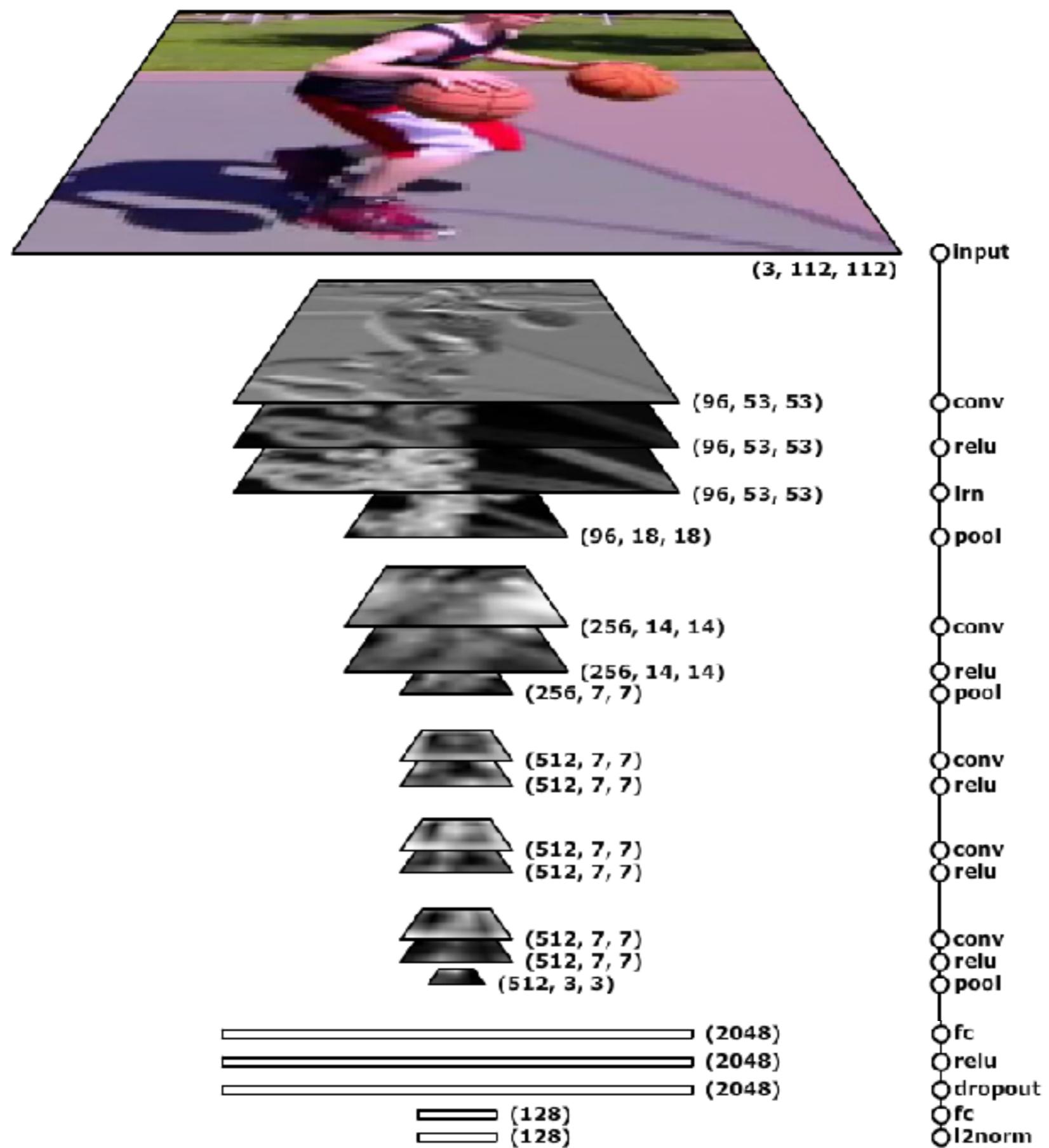
$$P(J_b^i; J_a^i)$$



$$\delta_{procrustes}(J_a, J_b) = \sum_{i=1}^N \sqrt{|J_a^i - P(J_b^i; J_a^i)|^2}$$







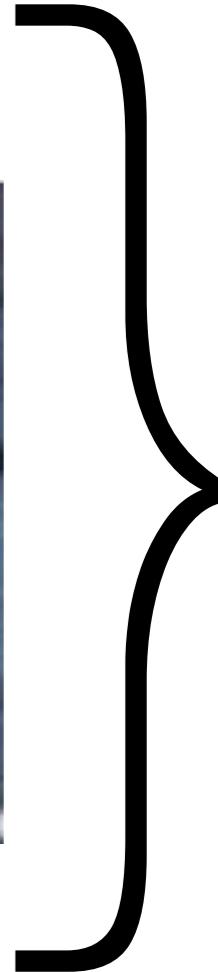
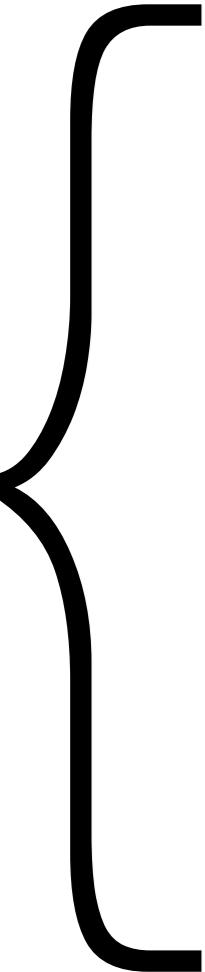
a

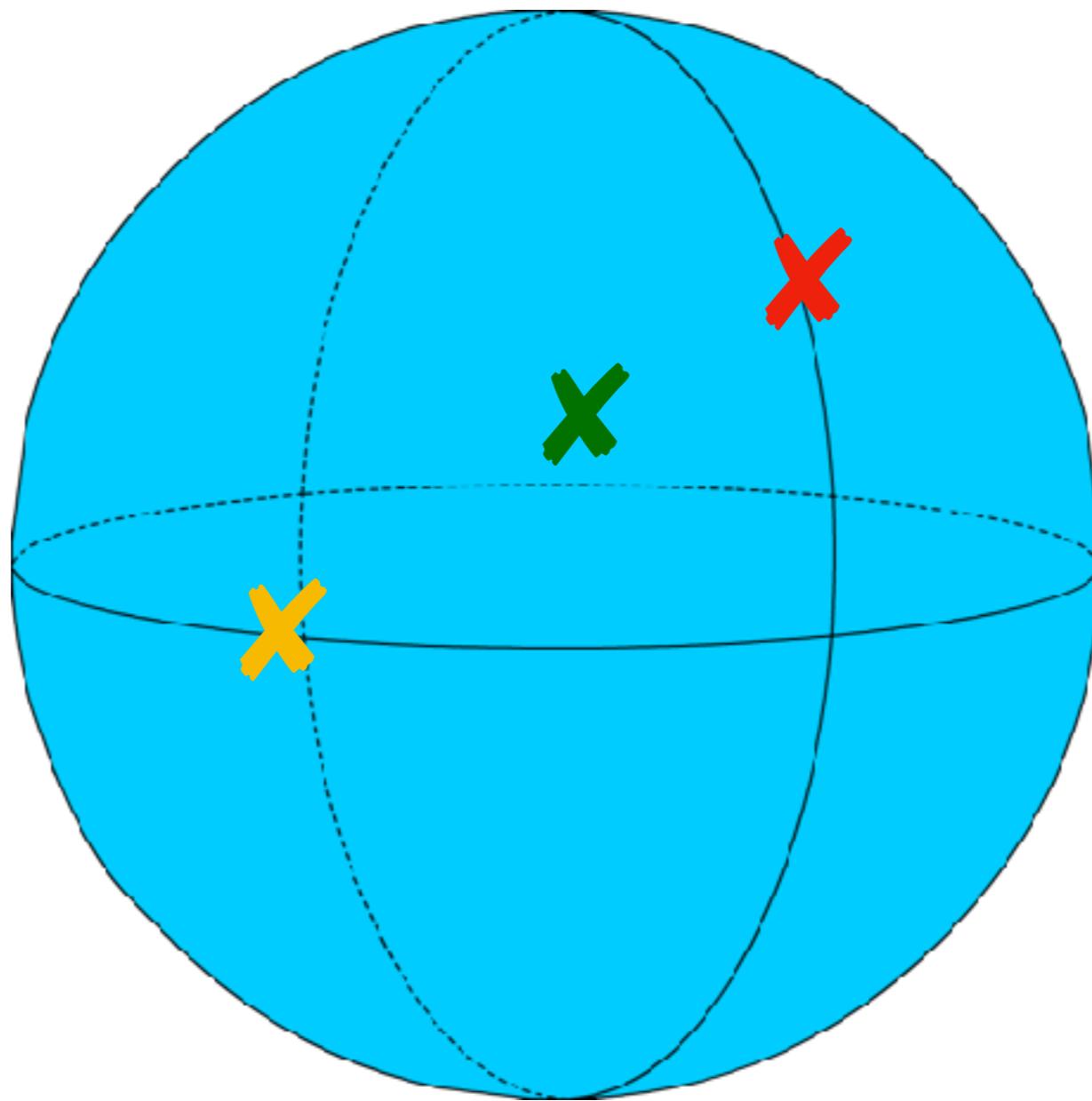


p

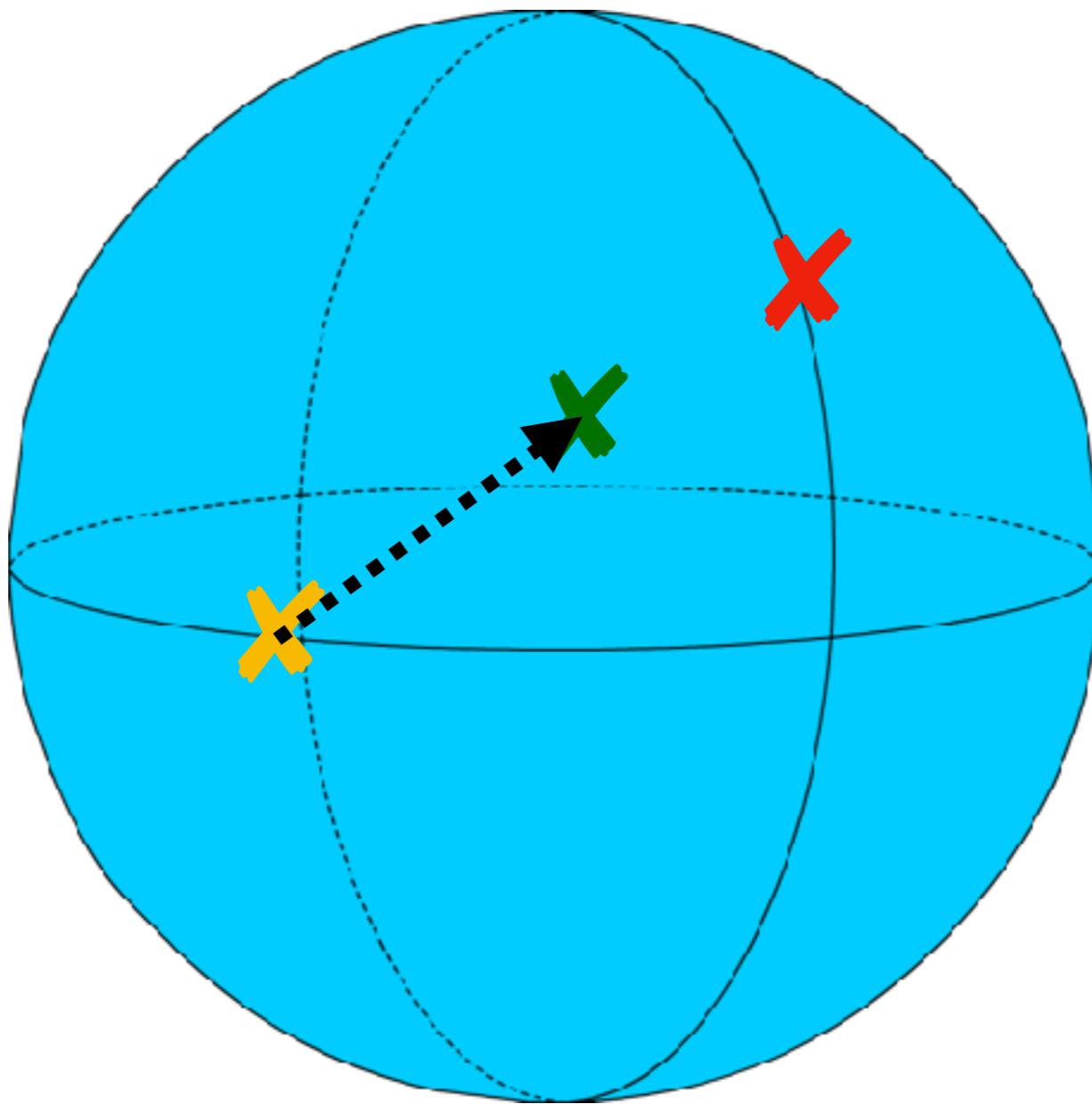


n

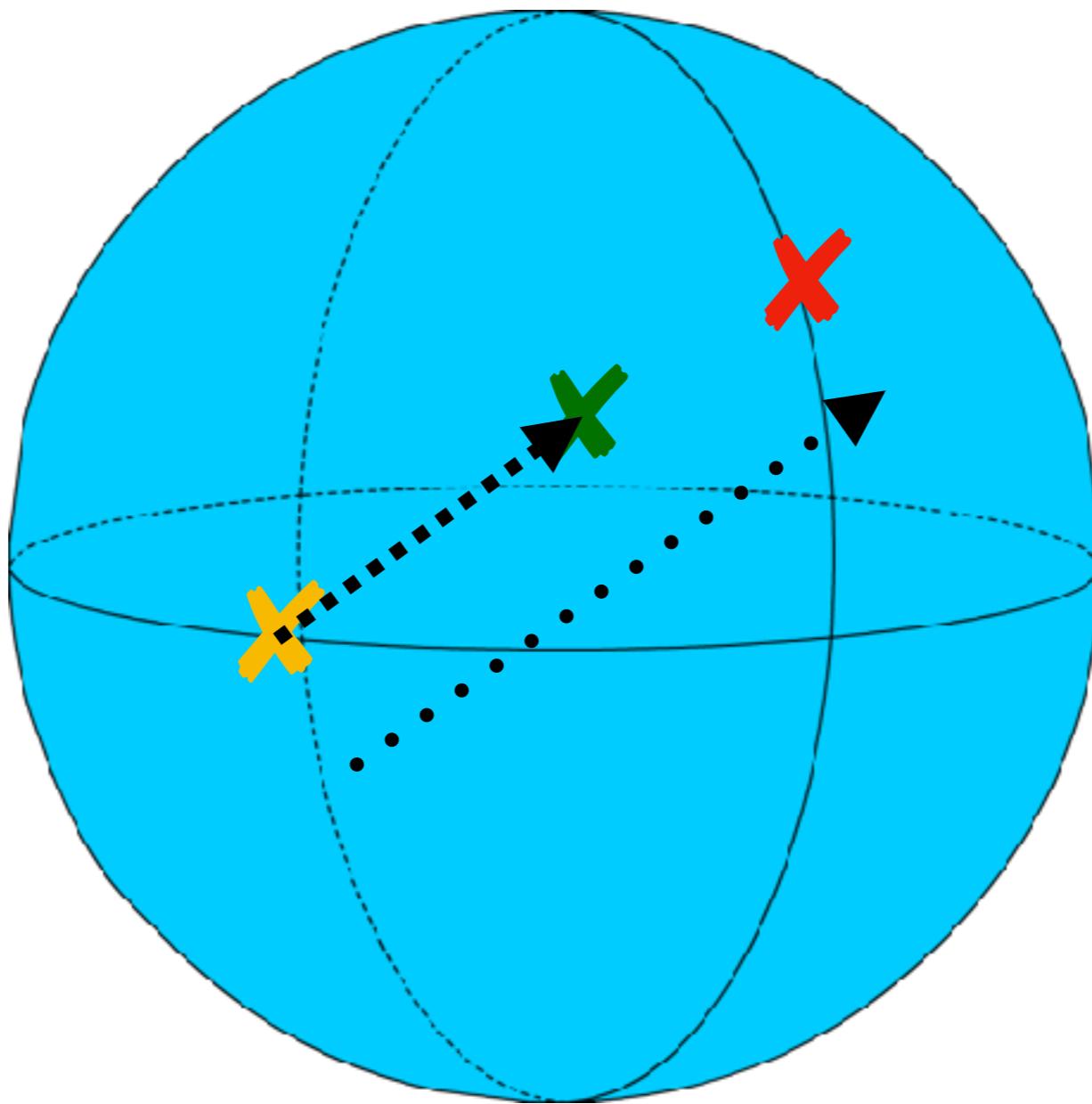




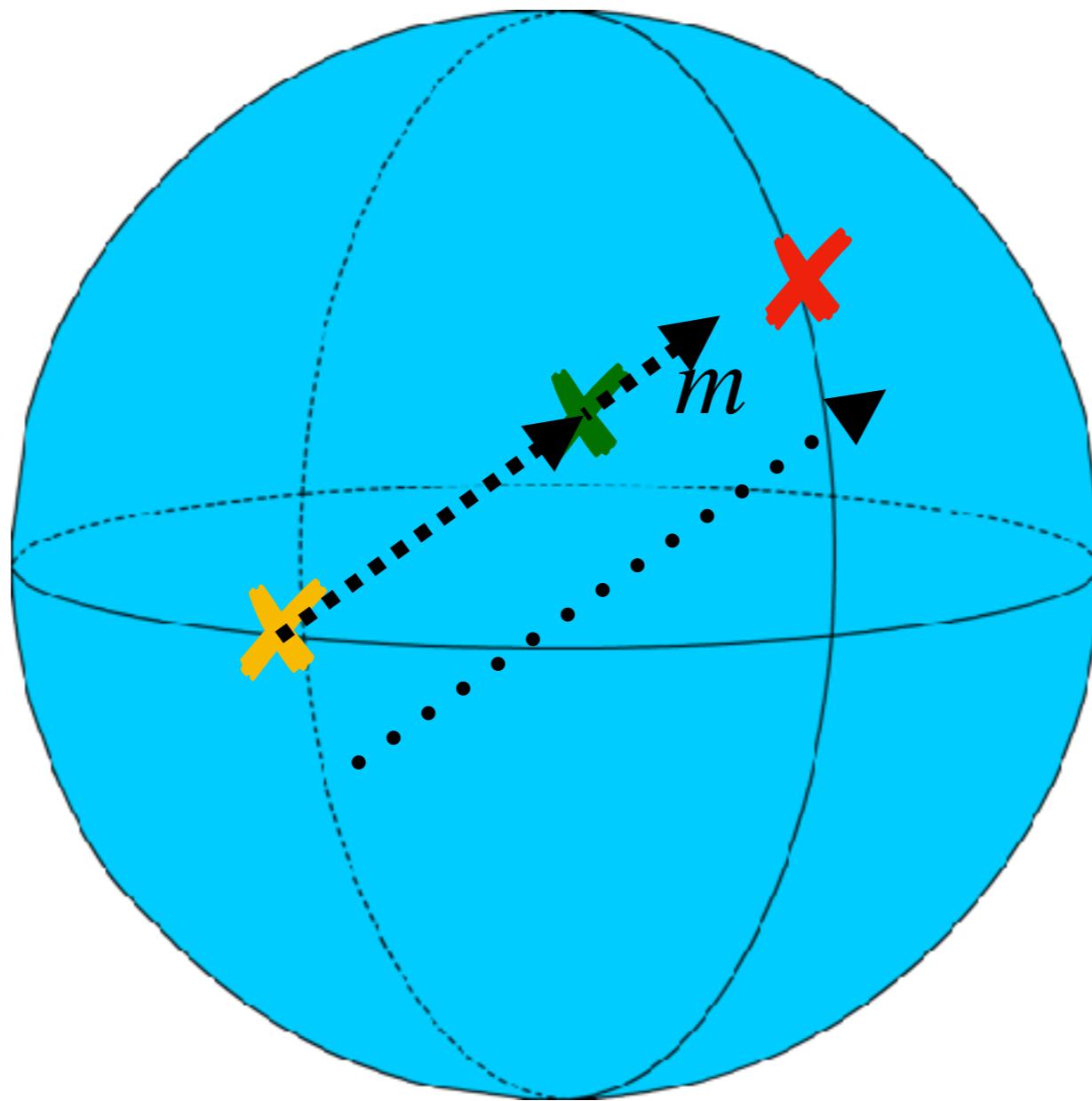
$$\mathcal{L}_{triplet} = [||f(a) - f(p)||_2^2 - ||f(a) - f(n)||_2^2 + m]_+$$



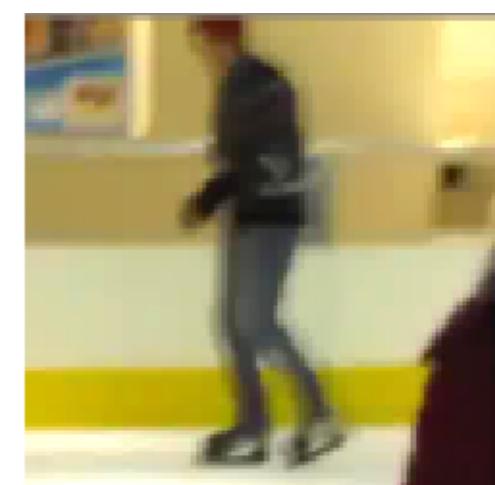
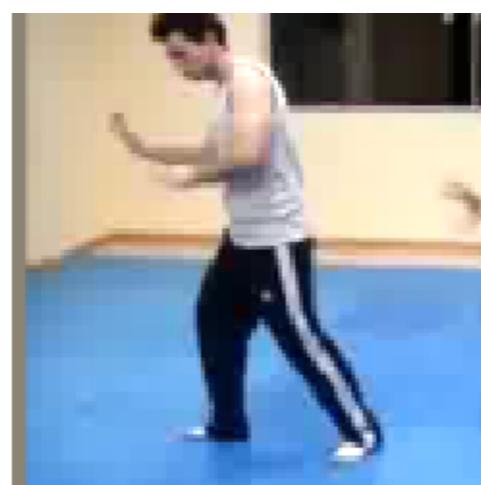
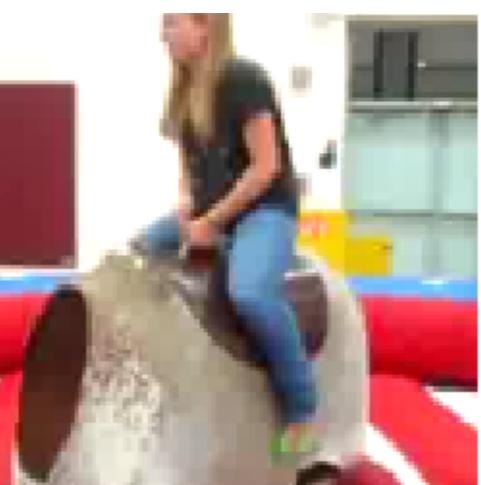
$$\mathcal{L}_{triplet} = [||f(a) - f(p)||_2^2 - ||f(a) - f(n)||_2^2 + m]_+$$

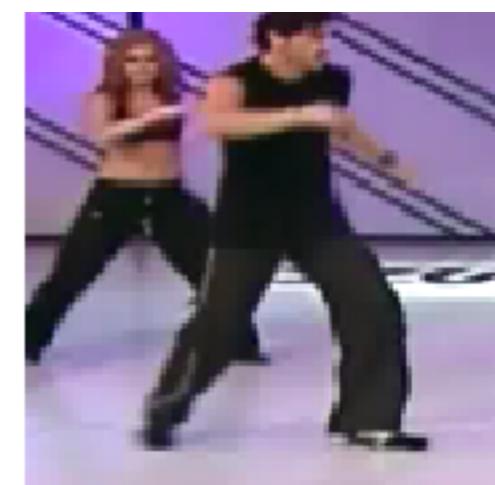
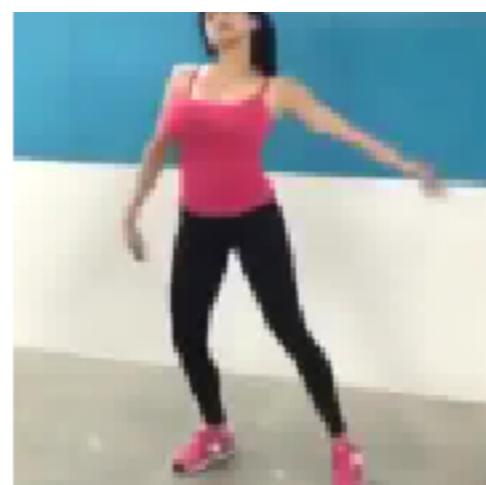
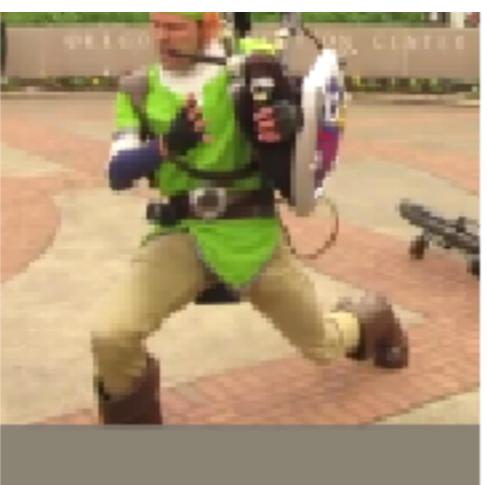
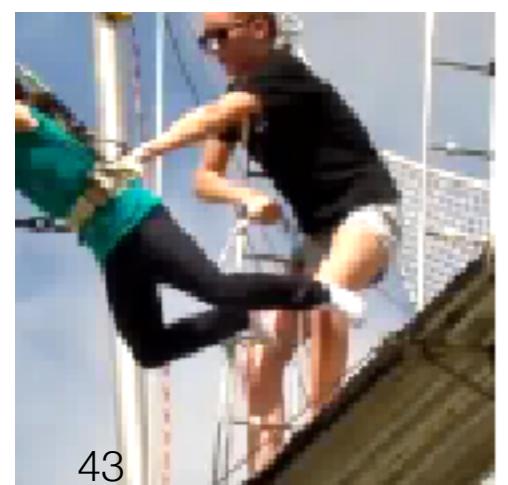
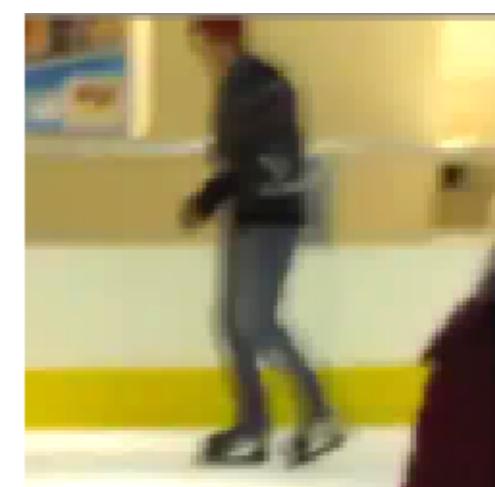
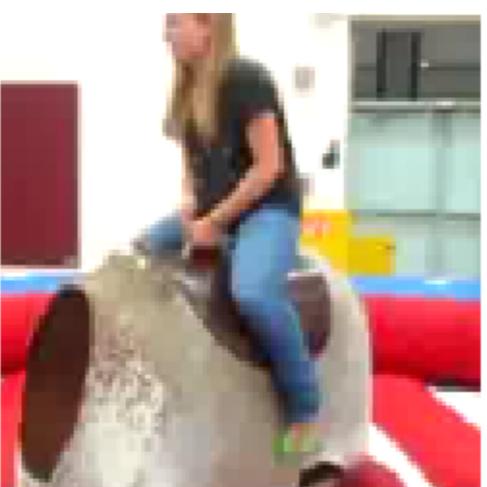


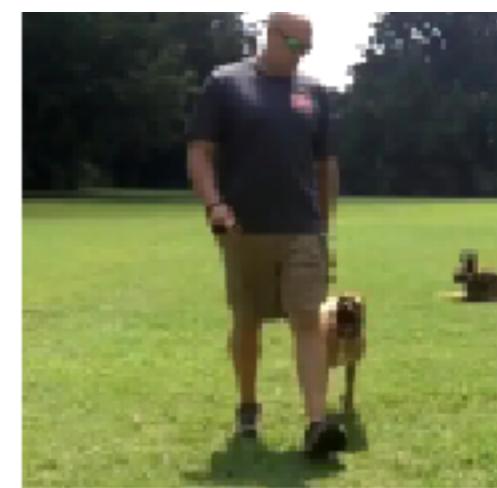
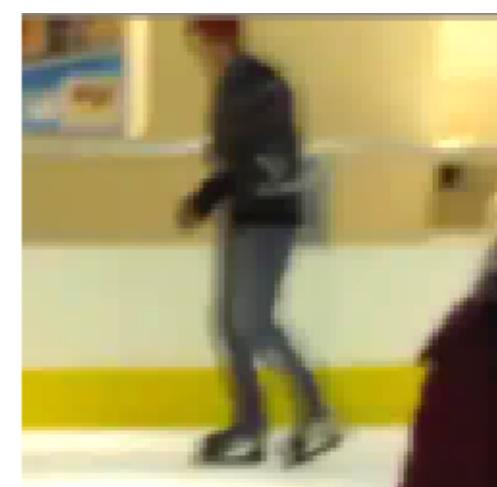
$$\mathcal{L}_{triplet} = [||f(a) - f(p)||_2^2 - ||f(a) - f(n)||_2^2 + m]_+$$

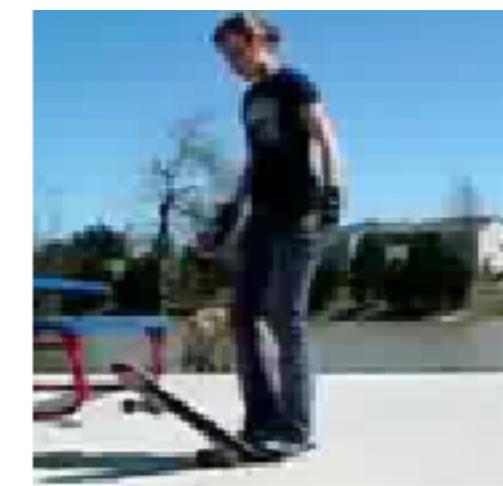
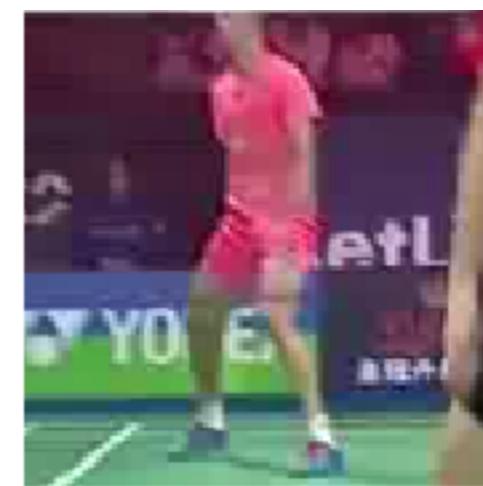
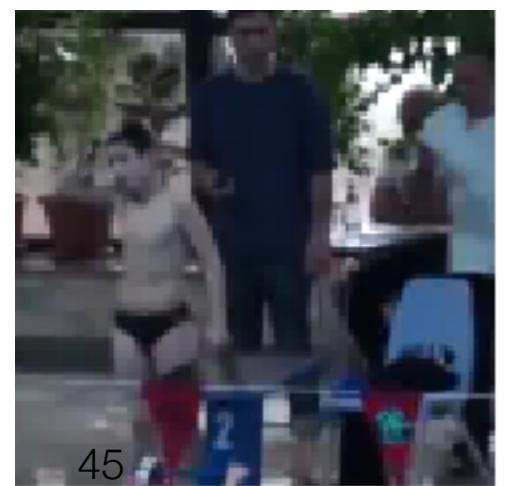
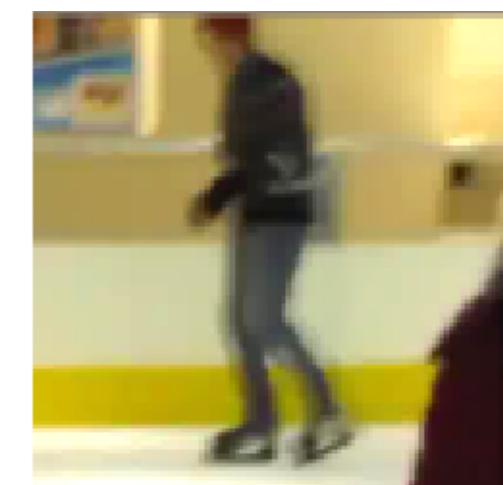
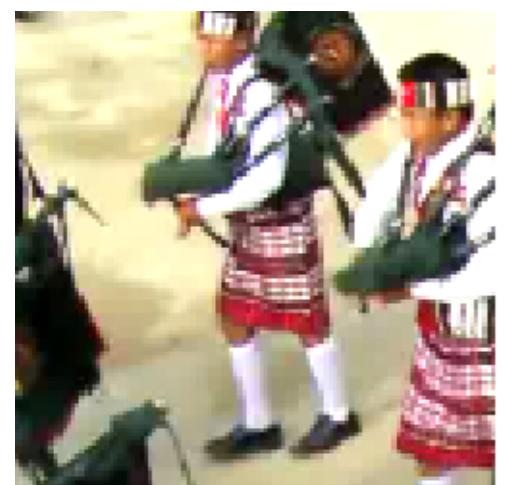


$$\mathcal{L}_{triplet} = [||f(a) - f(p)||_2^2 - ||f(a) - f(n)||_2^2 + m]_+$$





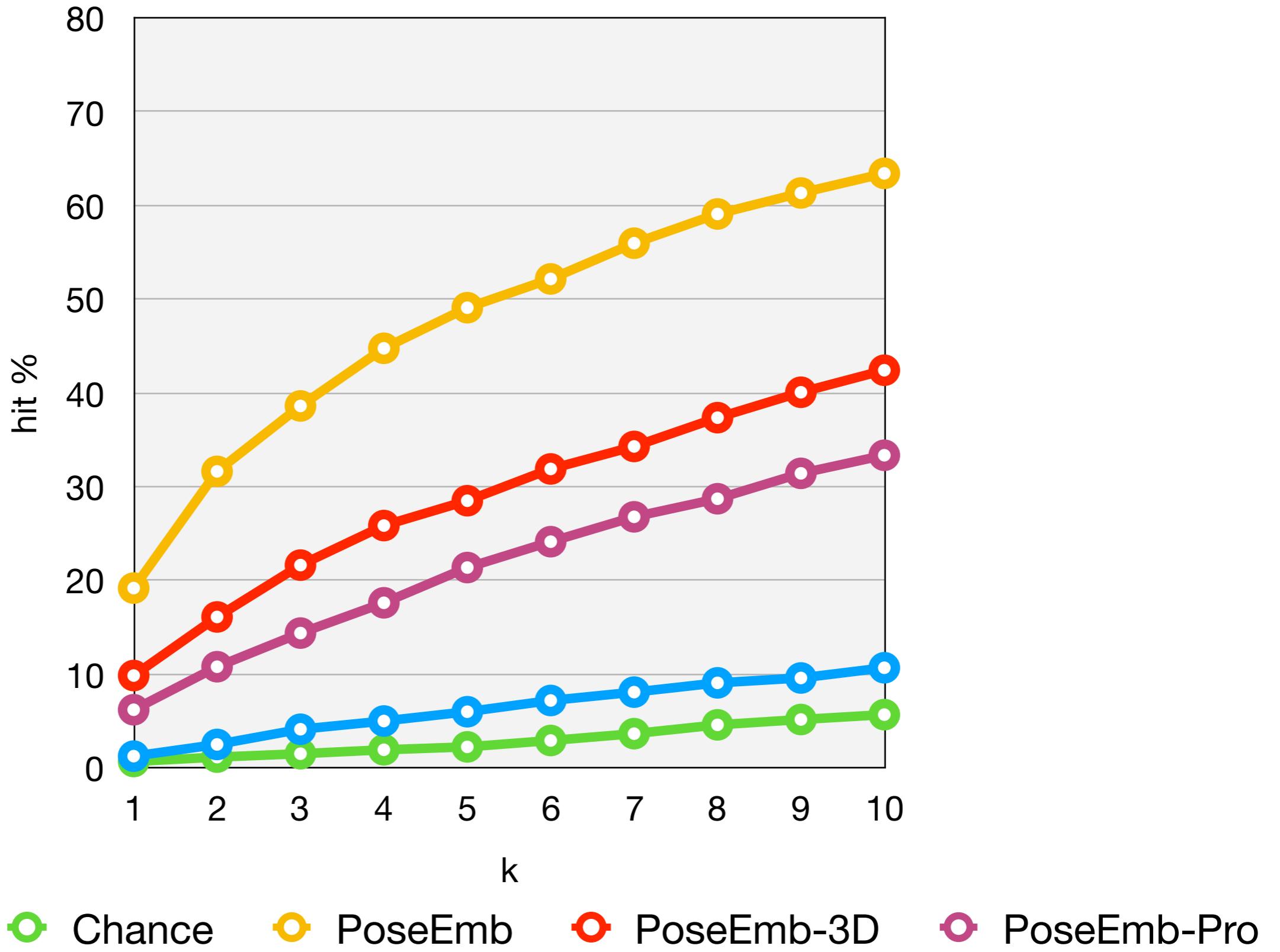




$$hit@k = \frac{1}{N} \sum_{i=1}^N hit_k(R_{1\dots k}^i)$$

$$hit_k(X) = \begin{cases} 1, & \text{if } |X \cap T_{50}| > 0 \\ 0, & \text{otherwise} \end{cases}$$

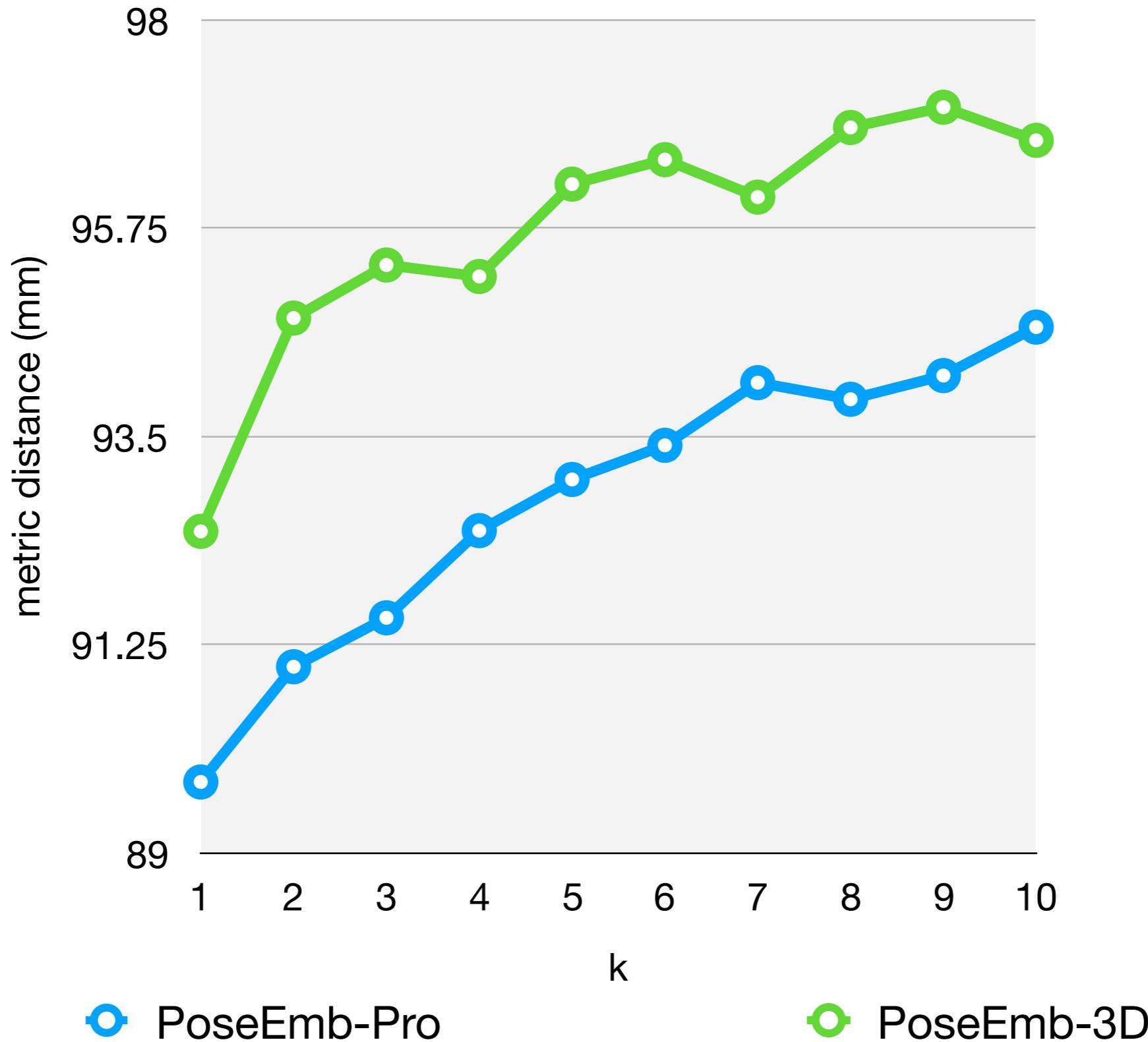
hit@k



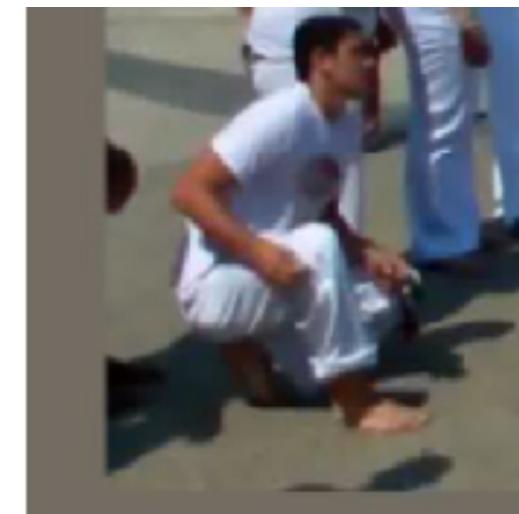
$$distance@k = \frac{1}{N} \sum_{i=1}^N distance_k(R_{1..k}^i)$$

$$distance_k(X) = \frac{1}{N} \sum_{i=1}^N \delta(X_i, q)$$

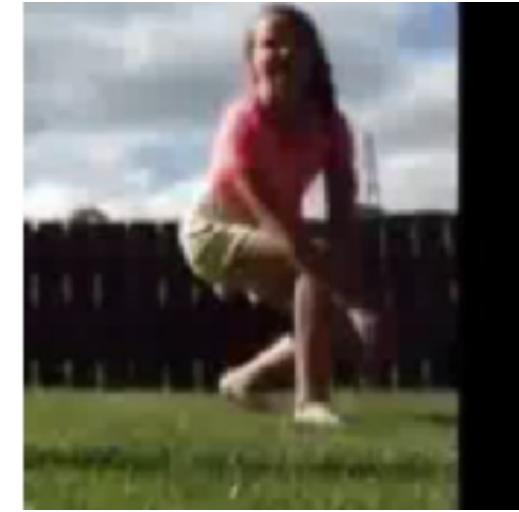
distance@k



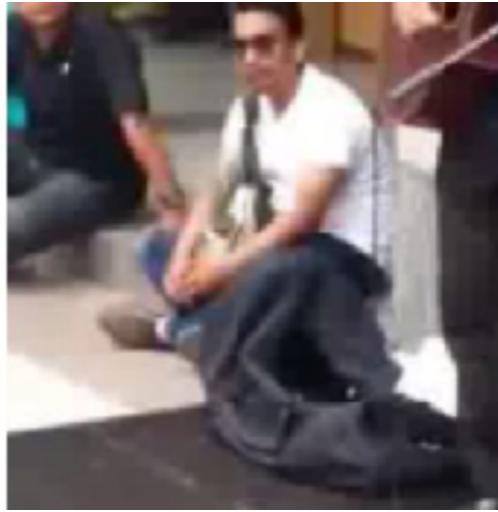
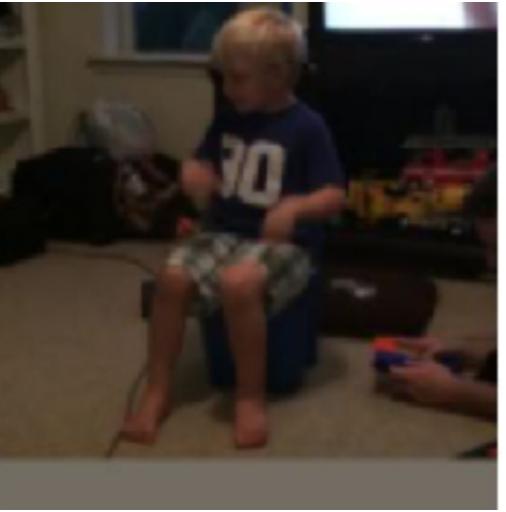
2D



3D



Procrustes



Language Querying

“Put your hands above your head!”



Posebits for Monocular Human Pose Estimation

Gerard Pons-Moll
MPI for Intelligent Systems
Tübingen, Germany
gerard.pons.moll@tue.mpg.de

David J. Fleet
University of Toronto
Toronto, Canada
fleet@cs.toronto.edu

Bodo Rosenhahn
Leibniz University of Hannover
Hannover, Germany
rosenhan@tnt.uni-hannover.de

Abstract

We advocate the inference of qualitative information about 3D human pose, called *posebits*, from images. *Posebits* represent boolean geometric relationships between body parts (e.g. left-leg in front of right-leg or hands close to each other). The advantages of *posebits* as a mid-level representation are 1) for many tasks of interest, such qualitative pose information may be sufficient (e.g. semantic image retrieval), 2) it is relatively easy to annotate large image corpora with *posebits*, as it simply requires answers to yes/no questions; and 3) they help resolve challenging pose ambiguities and therefore facilitate the difficult task of image-based 3D pose estimation. We introduce *posebits*, a *posebit* database, a method for selecting useful *posebits* for pose estimation and a structural SVM model for *posebit* inference. Experiments show the use of *posebits* for semantic image retrieval and for improving 3D pose estimation.

1. Introduction

While tremendous effort has focused on the extraction of quantitative 3D human pose from images and video, in this paper we consider the estimation of qualitative pose information, called *posebits*. *Posebits* are attributes of pose that specify the relative positions or orientations of body parts. They have the advantage that one can easily collect *posebit* image annotations for training purposes, and they can be reliably inferred from images. Further, they are useful for resolving 3D pose ambiguities and for myriad other tasks where quantitative pose is not required.

The effective use of both generative and discriminative approaches to pose estimation require training data com-

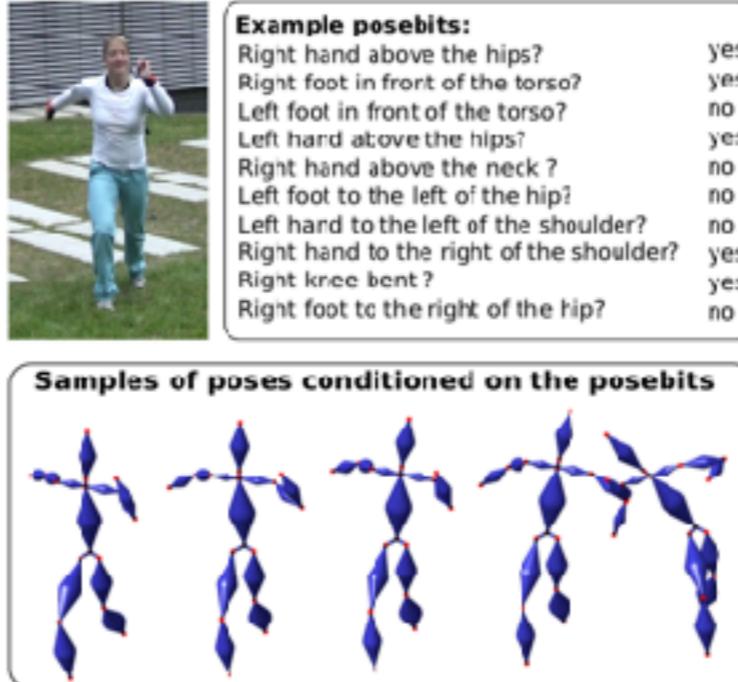
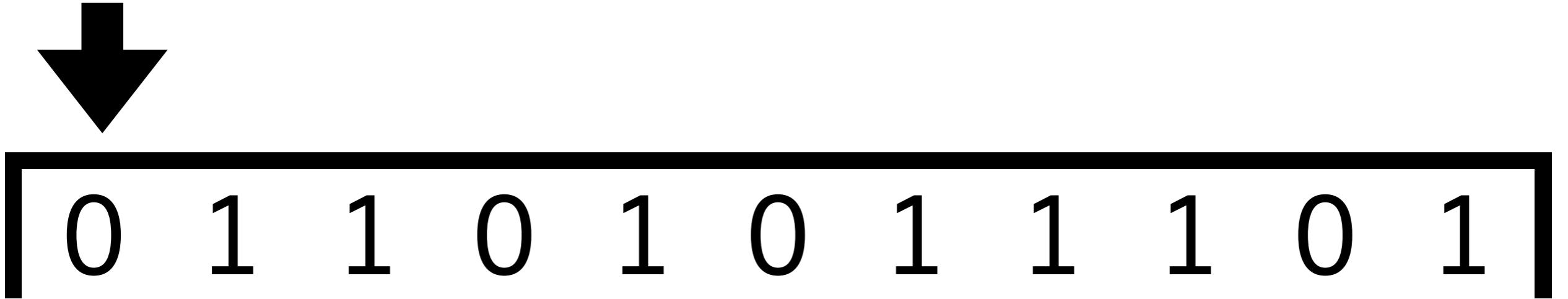


Figure 1. Top: *Posebits* are inferred directly from image features using a trained classifier. Since *posebits* consist of simple yes/no questions, images can be easily annotated by humans. They may be useful for many tasks. For example, on the bottom image, we show samples of poses conditioned on the *posebits* depicted on the top image. By conditioning the poses on *posebits* uncertainty about the pose is reduced. Notice how the poses are qualitatively very similar to the observed image. In this example we show the ground truth posebyte. Our model also takes into account uncertainty in the estimation of *posebits* by marginalizing over them.

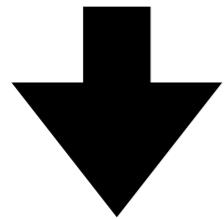
to errors [5, 6].

By contrast, it is relatively simple to obtain training data for *posebits* from human annotations. Indeed, people often

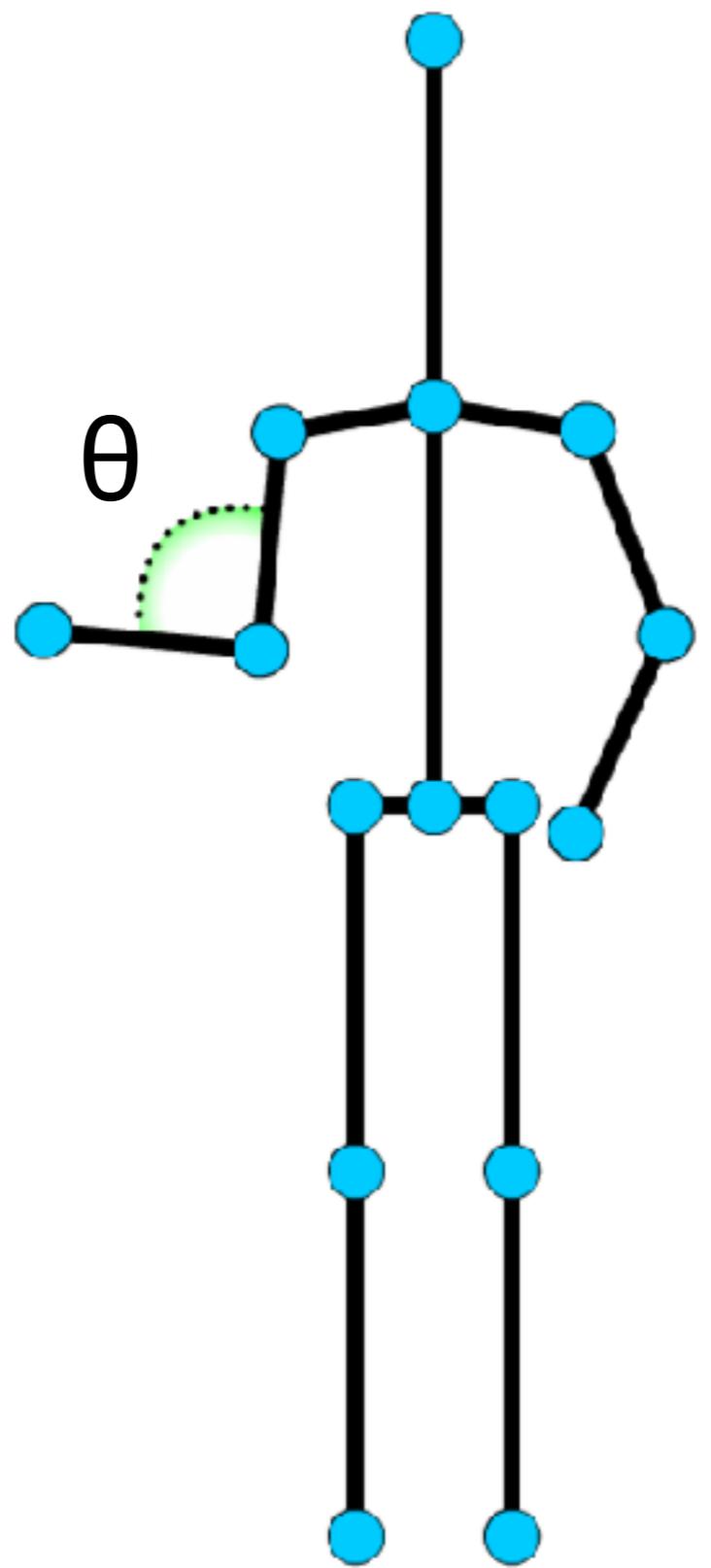
“Right knee **is not** bent”

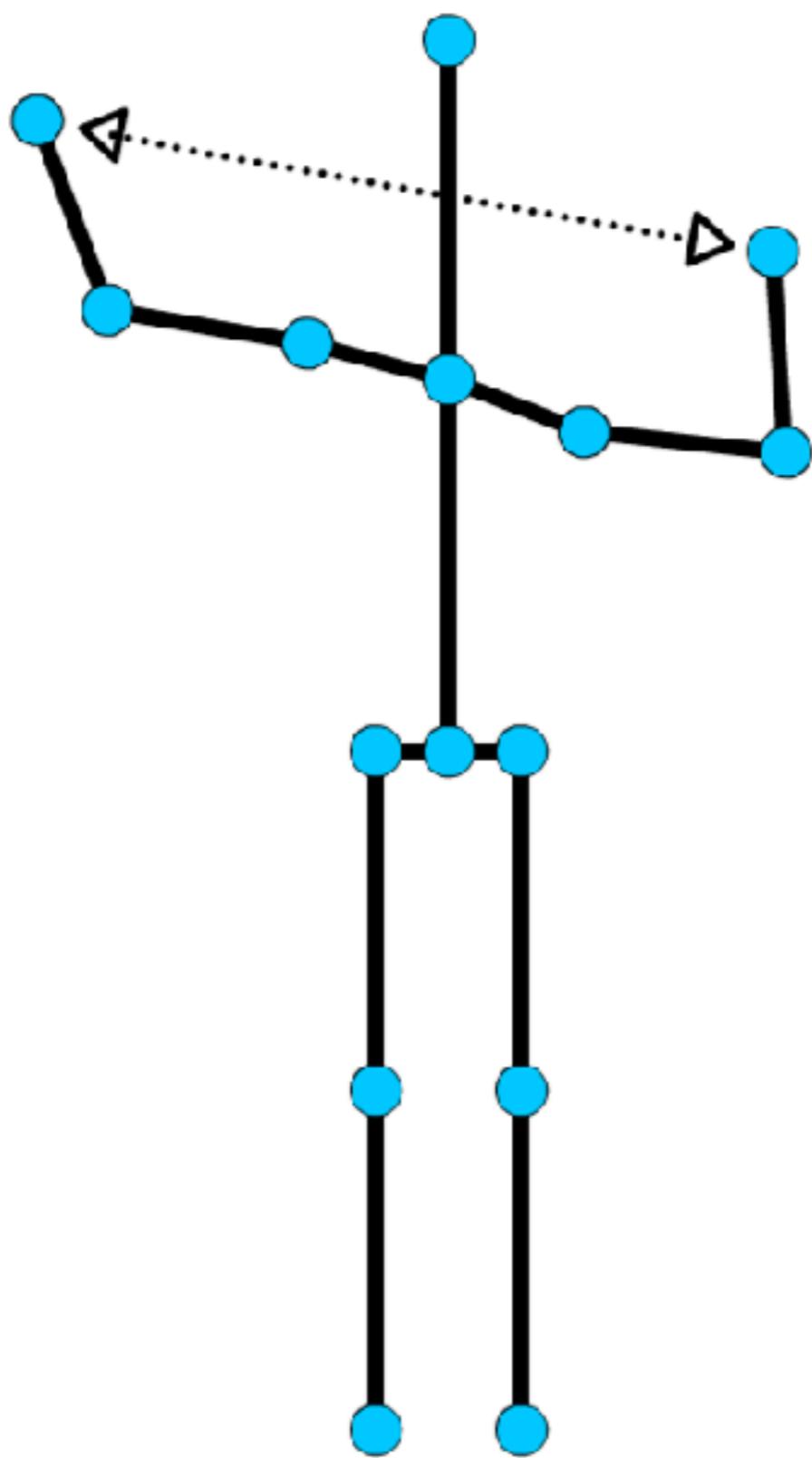


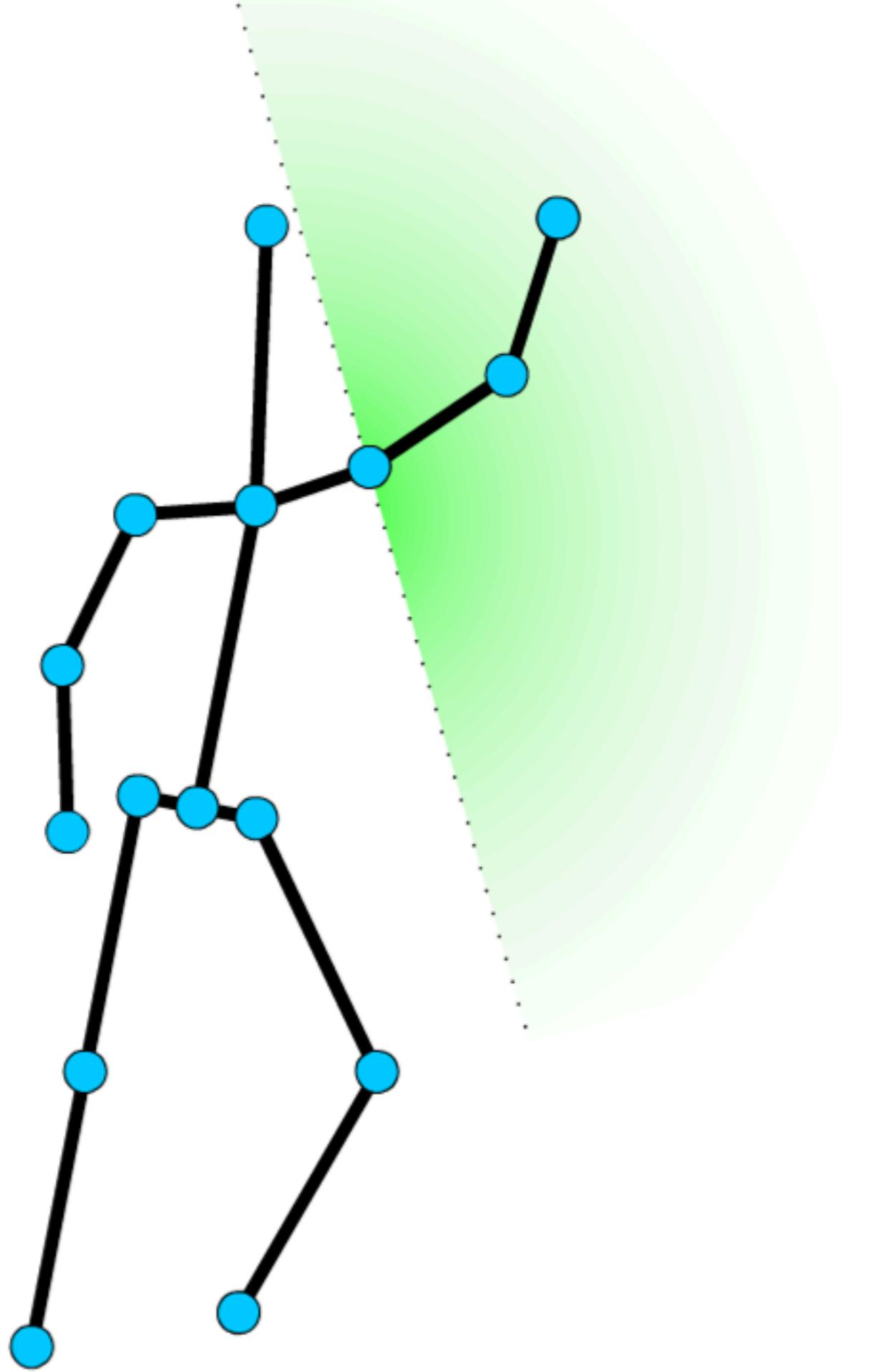
“Left elbow **is** bent”



0	1	1	0	1	0	1	1	1	0	1
---	---	---	---	---	---	---	---	---	---	---





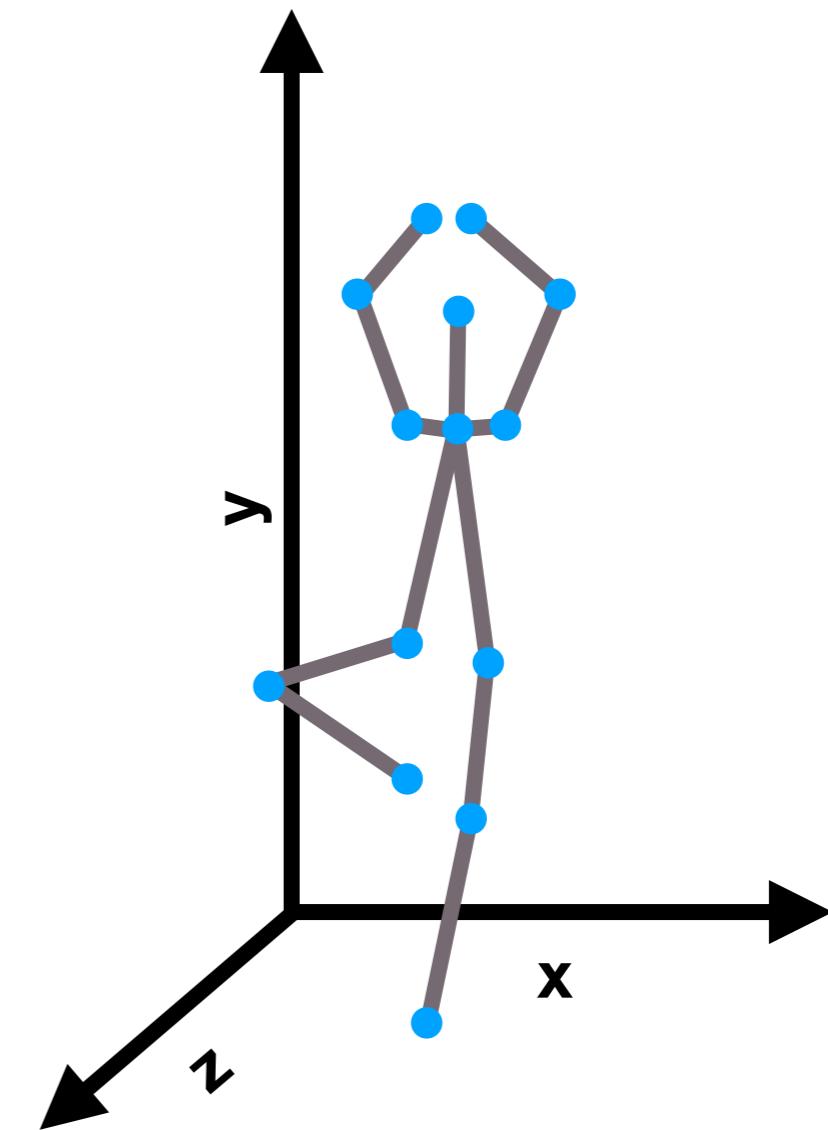
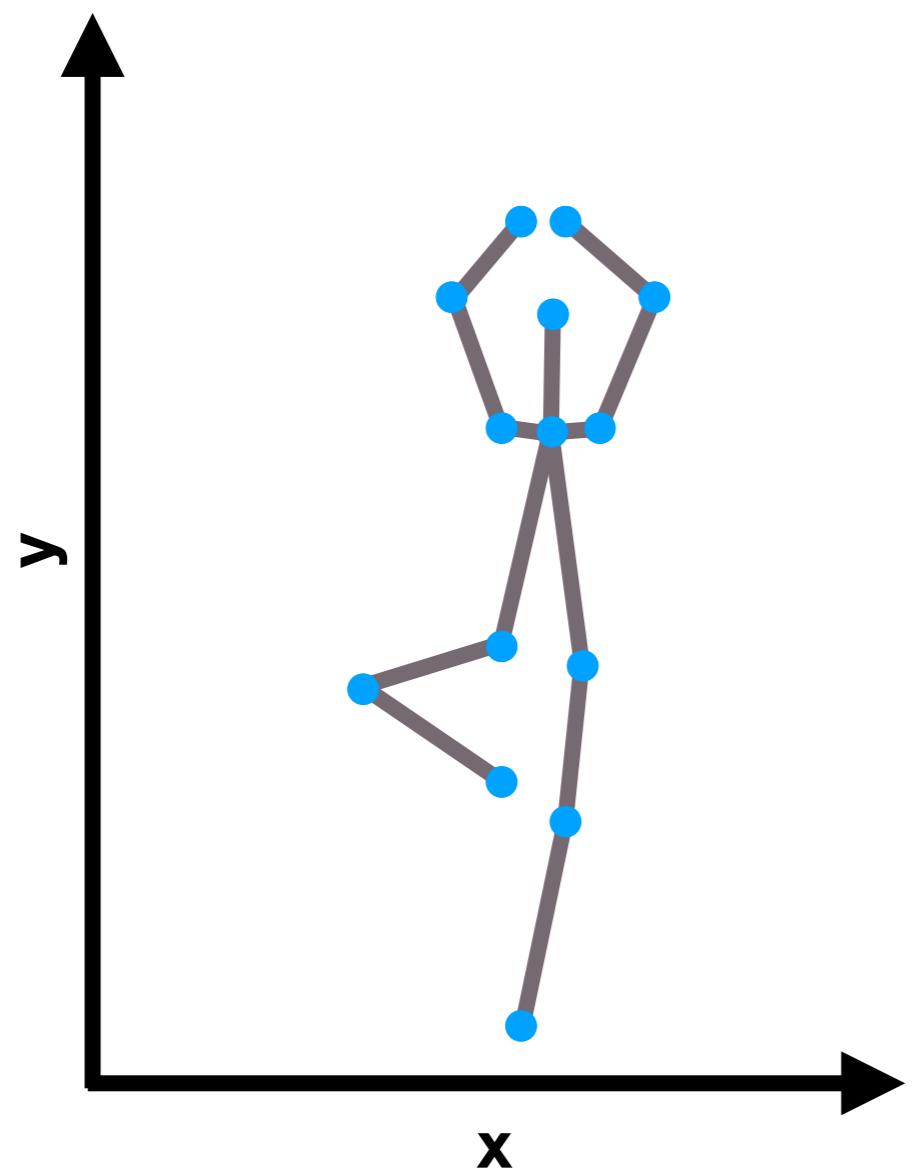


“Put your hands above your head!”

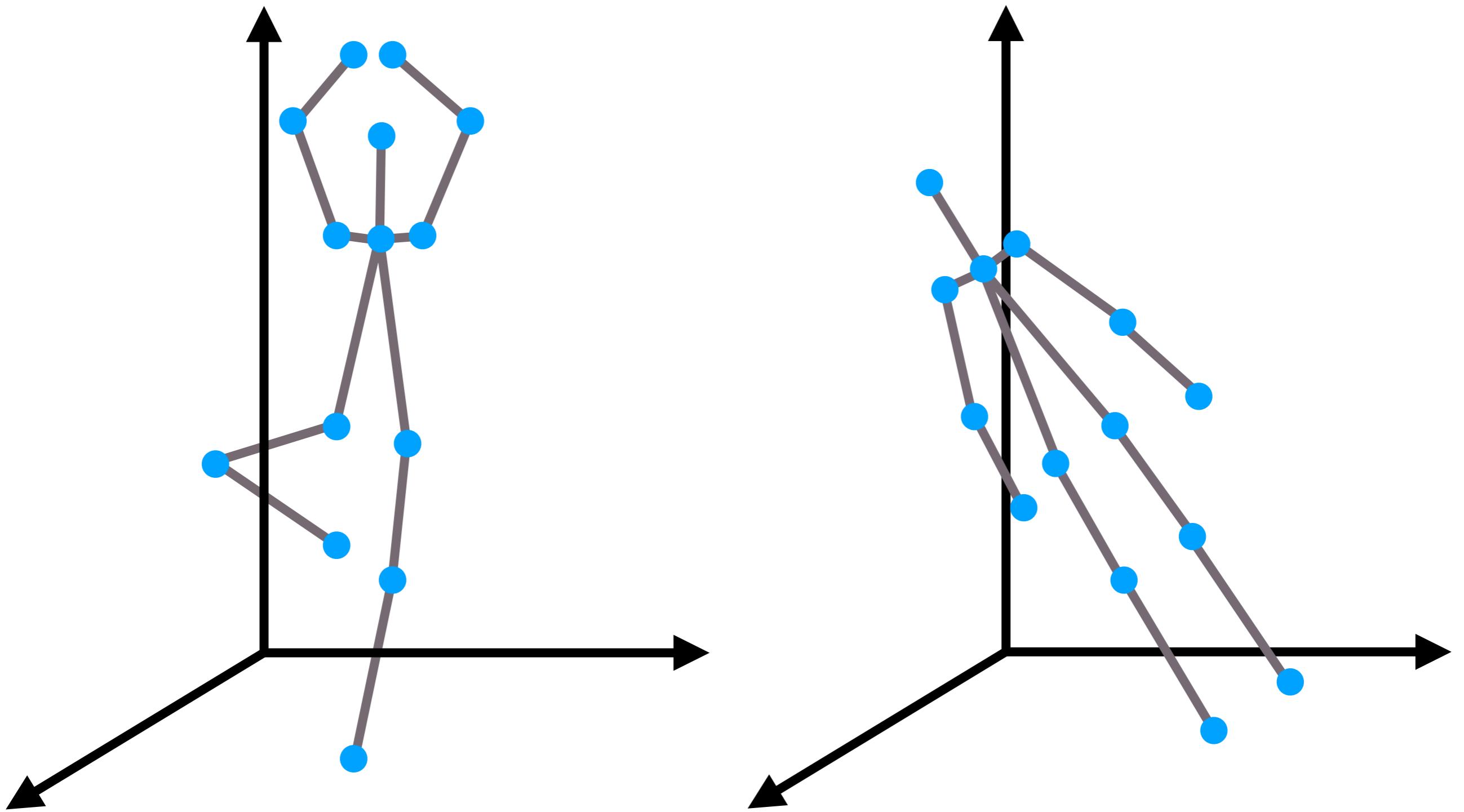


“Put your left hand above your head;
Put your right hand above your head!”





⋮ [1 0 0 1 0 1 1 1 0 1 0 1 1] ⋮



$$\delta_{hamming}(J_a, J_b) = \frac{1}{P} \sum J_a \oplus J_b$$

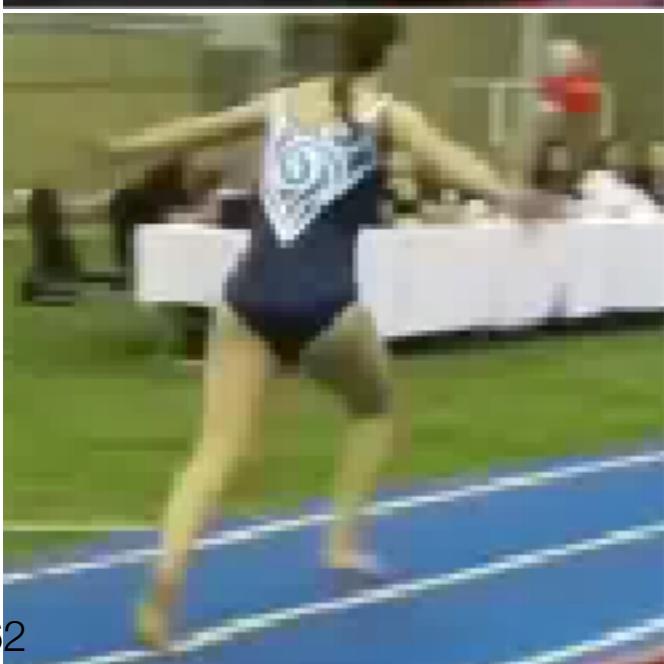
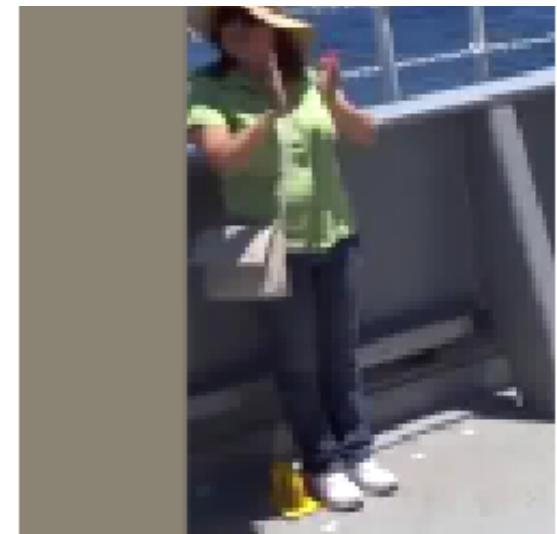
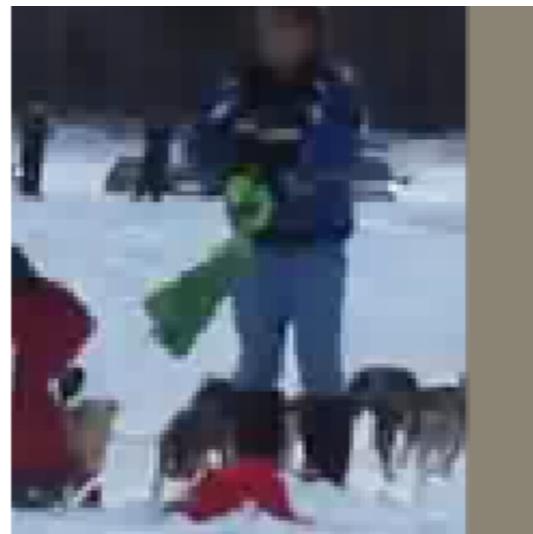
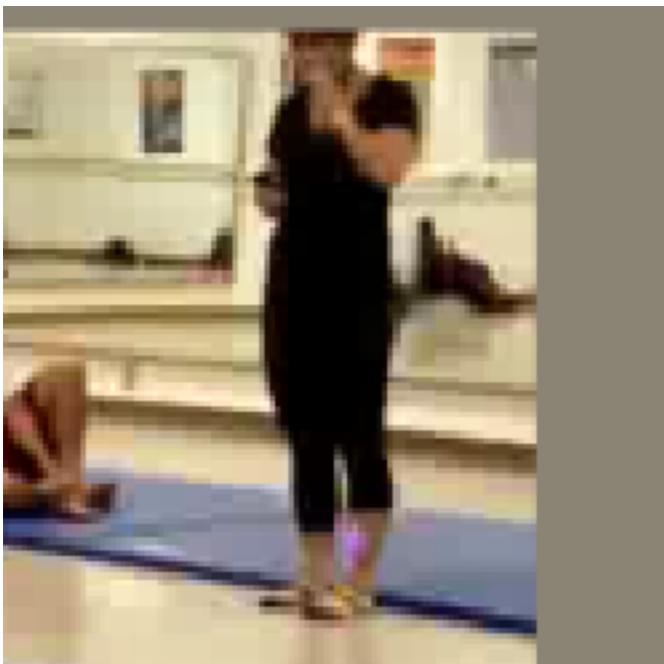
1
1
1
:
1
0
1
1

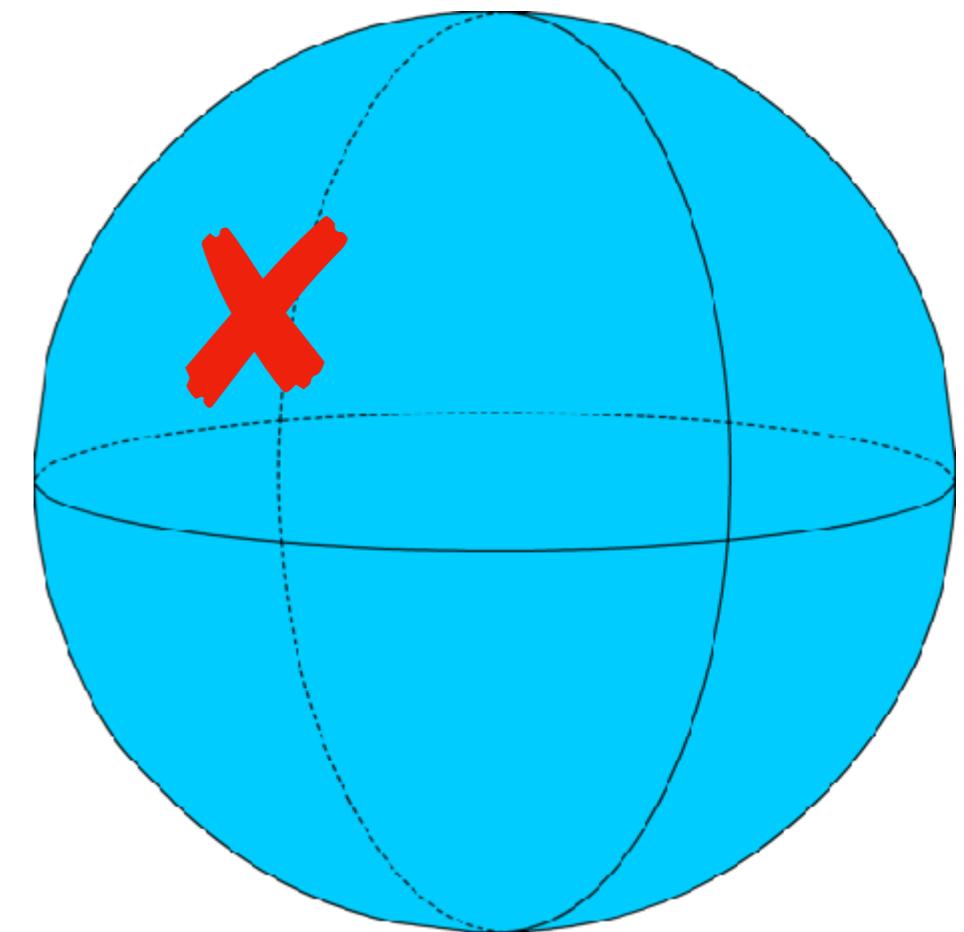
0
1
0
:
1
1
0
1

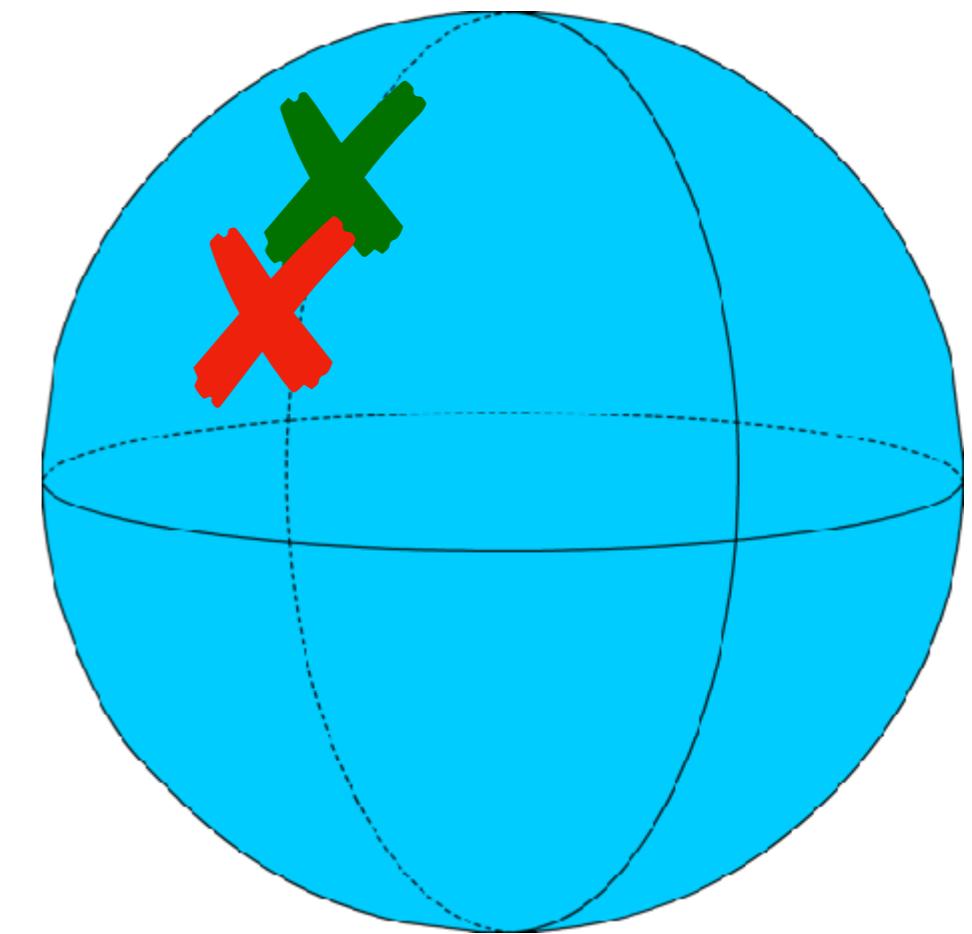
$$\delta_{hamming}(J_a, J_b) = \frac{1}{P} \sum J_a \oplus J_b$$

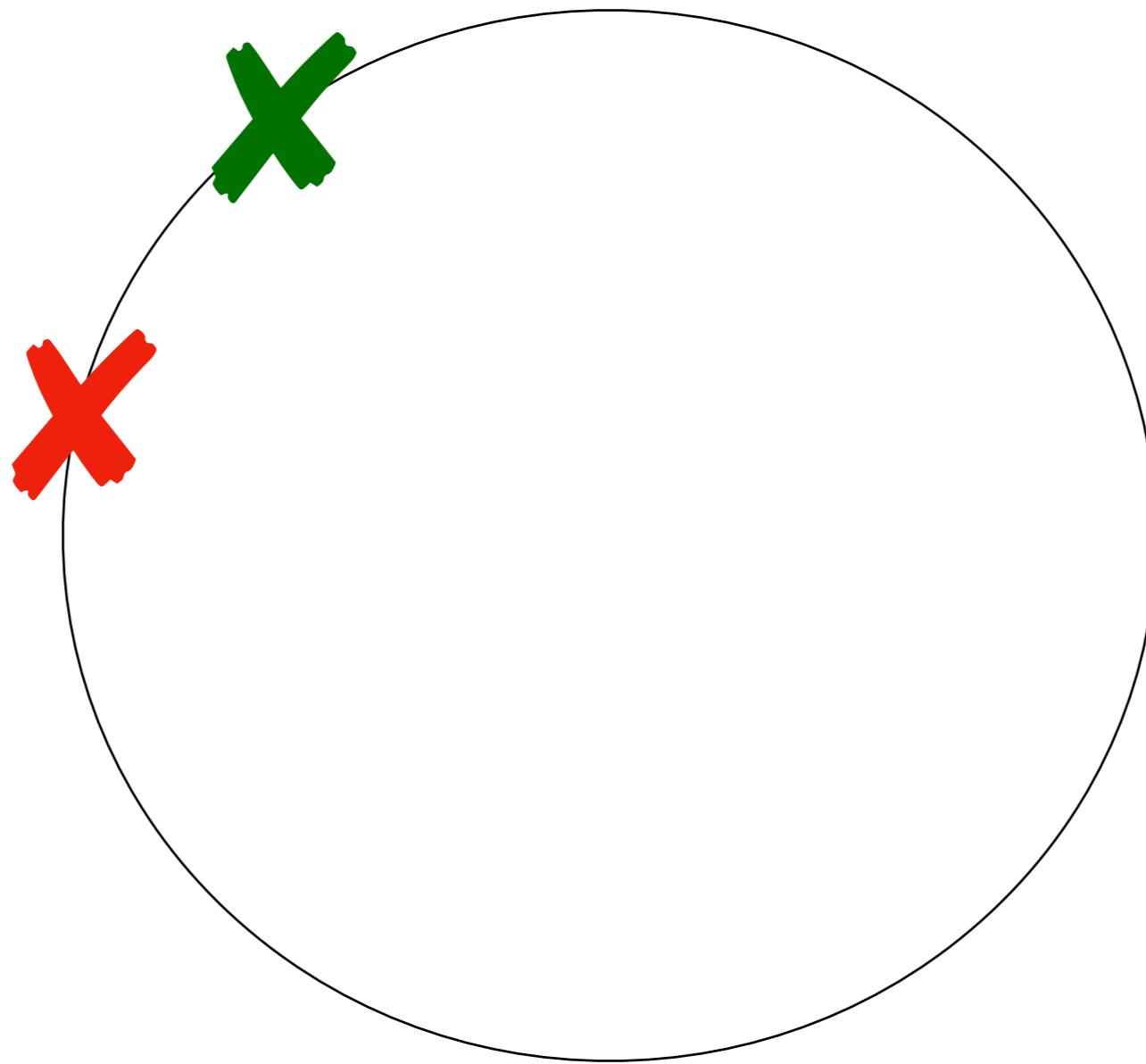
1⊕0
1⊕1
1⊕0
⋮
1⊕1
0⊕1
1⊕0
1⊕1

$$\delta_{hamming}(J_a, J_b) = \frac{1}{P} \sum J_a \oplus J_b$$

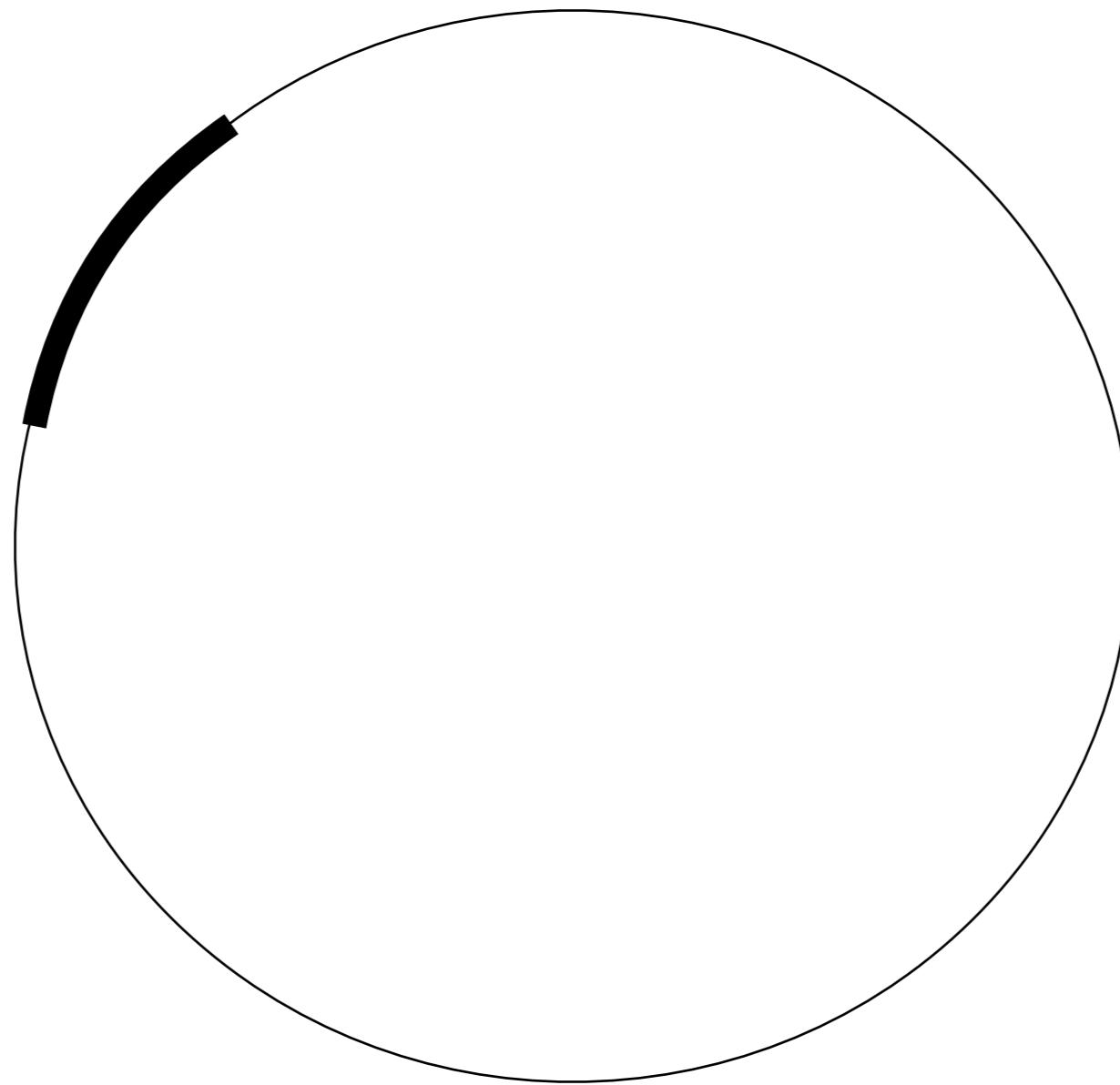




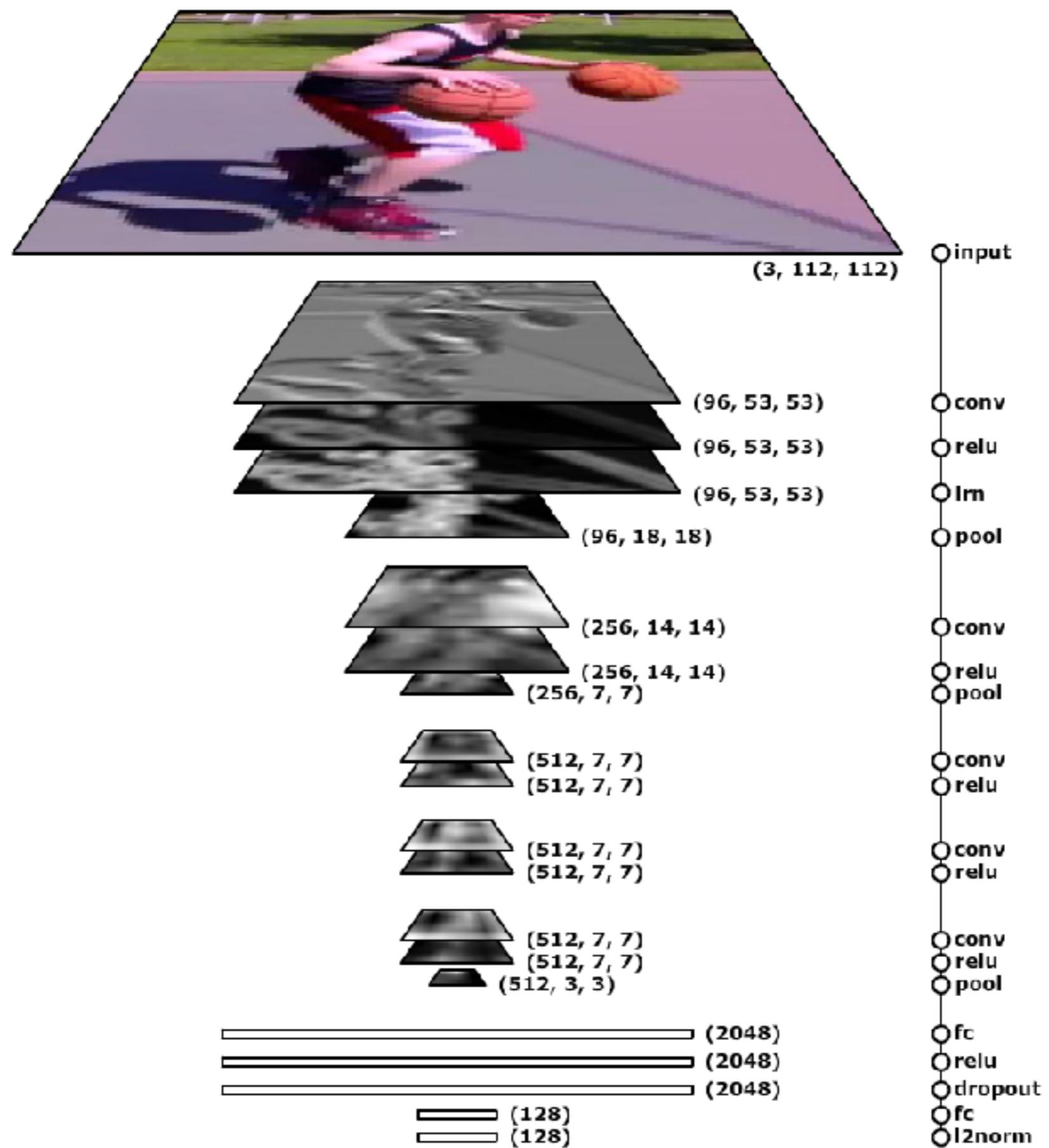
$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 1 \\ 1 \\ 0 \\ 1 \end{bmatrix}$$


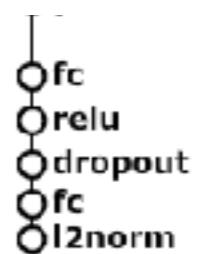
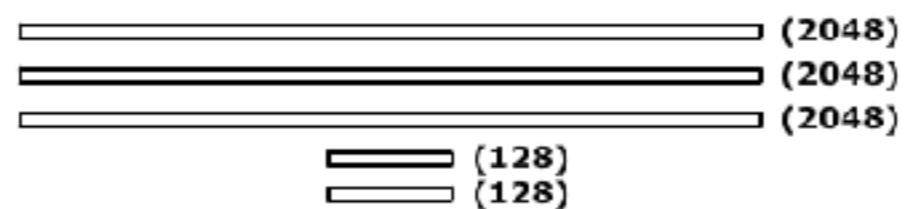


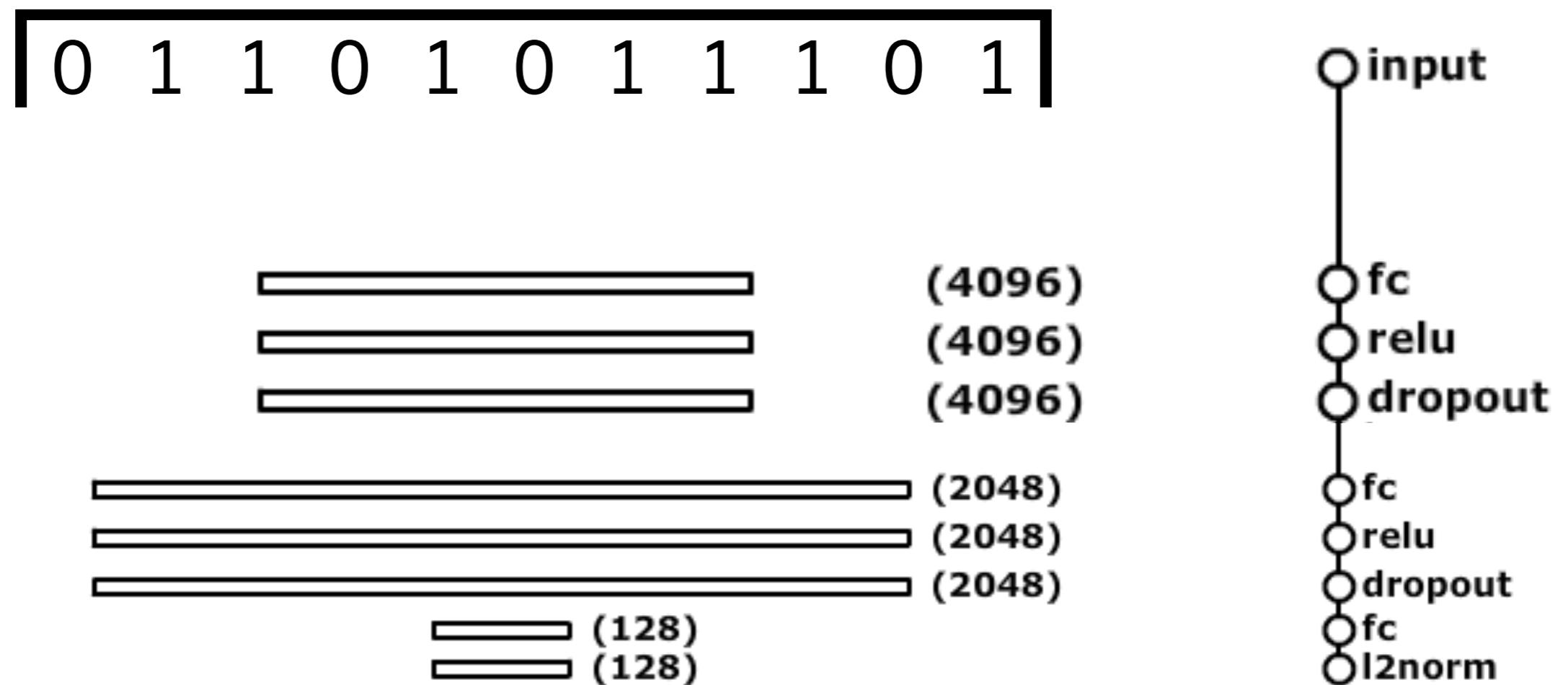
$$\mathcal{L}_c = 1 - f(e_{image}) \cdot g(e_{pose})$$

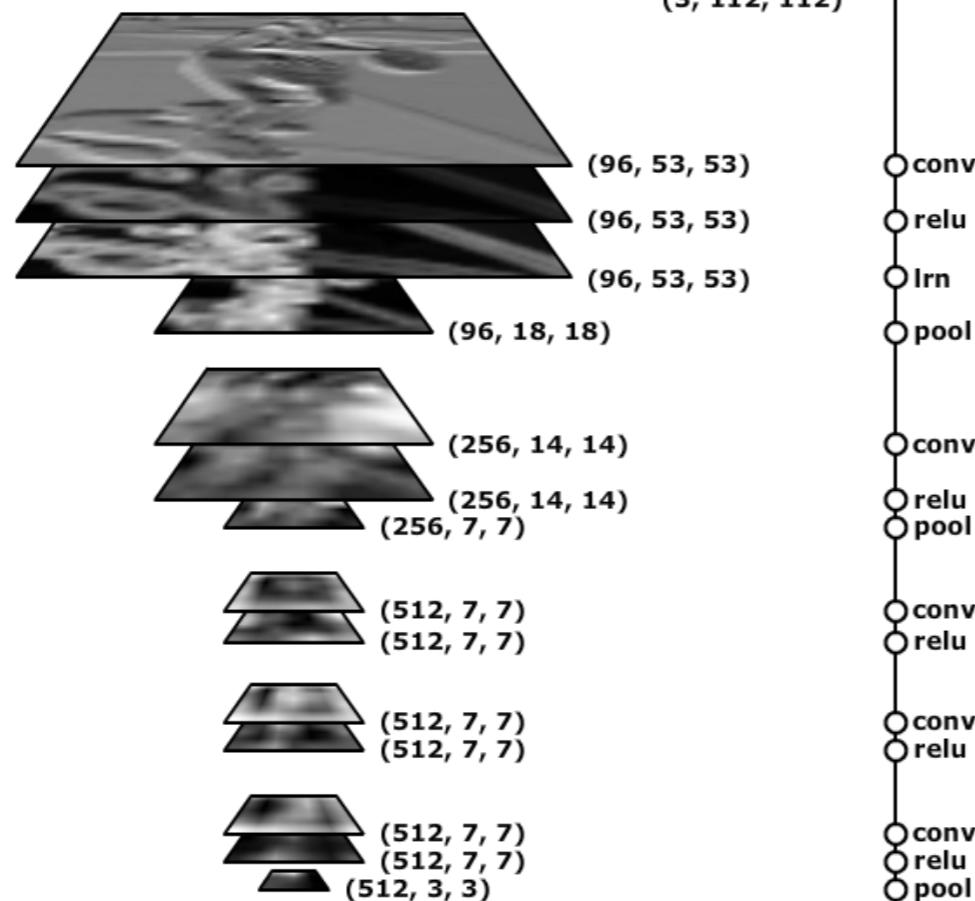


$$\mathcal{L}_c = 1 - f(e_{image}) \cdot g(e_{pose})$$

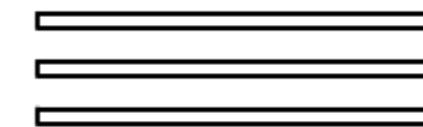






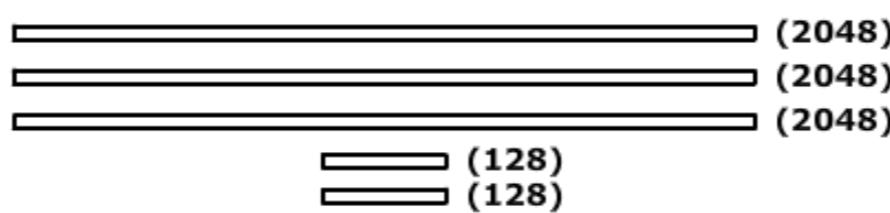


1 0 1 ... 0 1 1 0



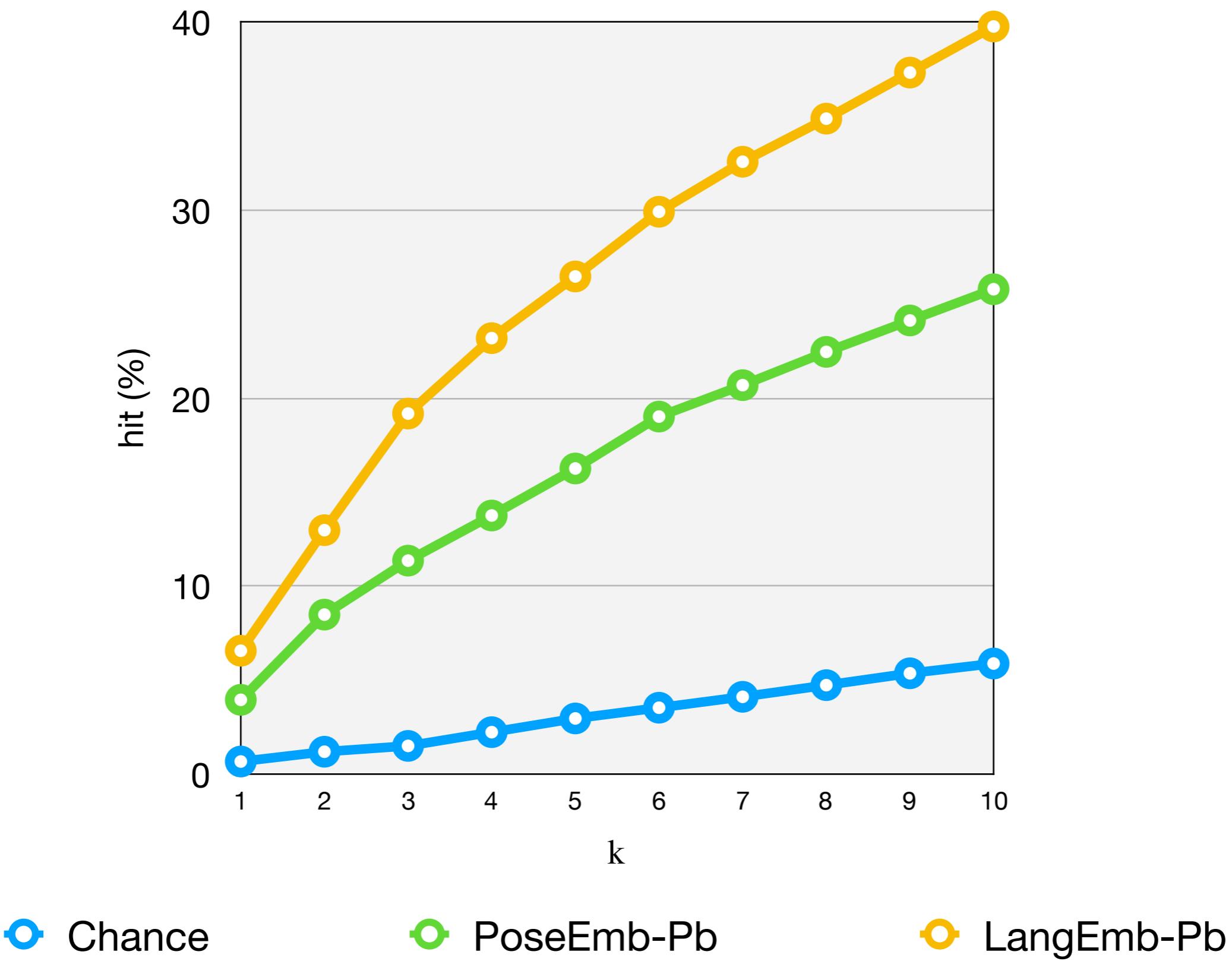
(4096)
(4096)
(4096)

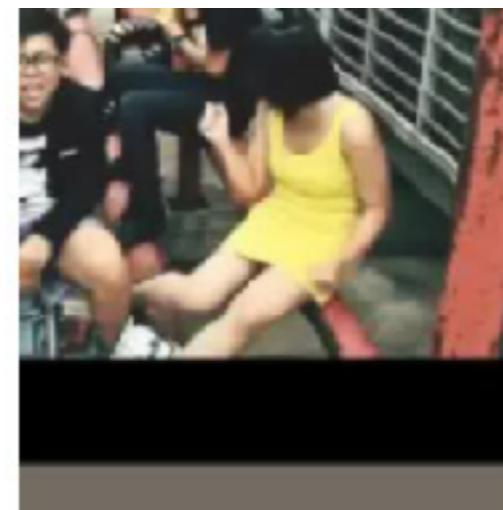
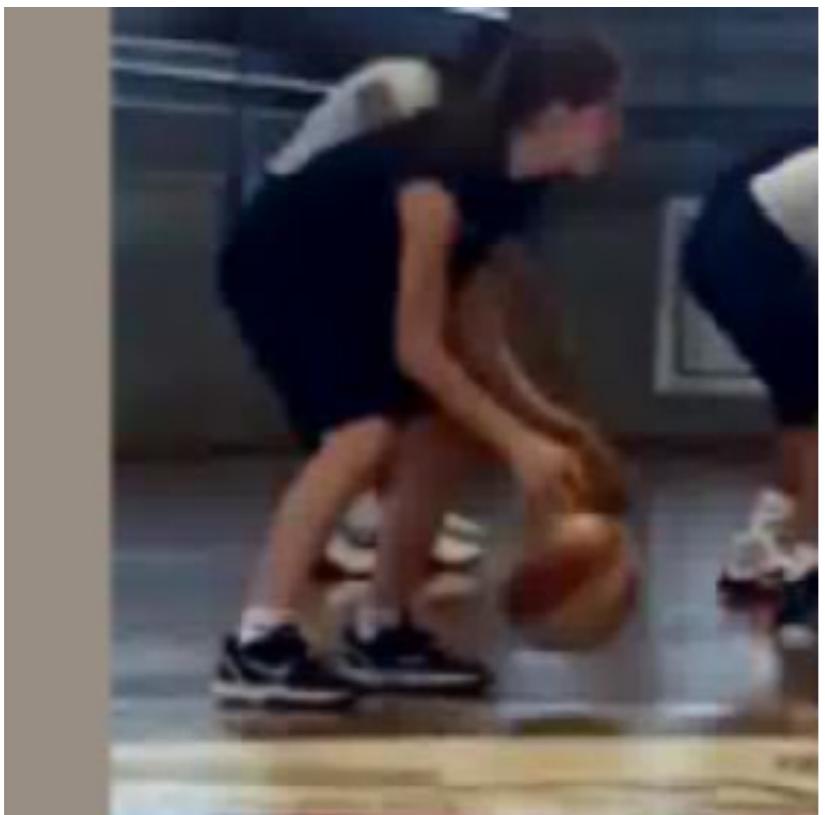
input
fc
relu
dropout



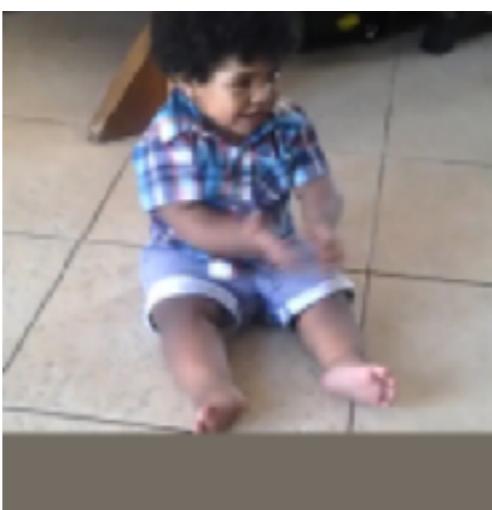
fc
relu
dropout
fc
l2norm

hit@k





1
0
0
:
1
1



Language Subset Querying



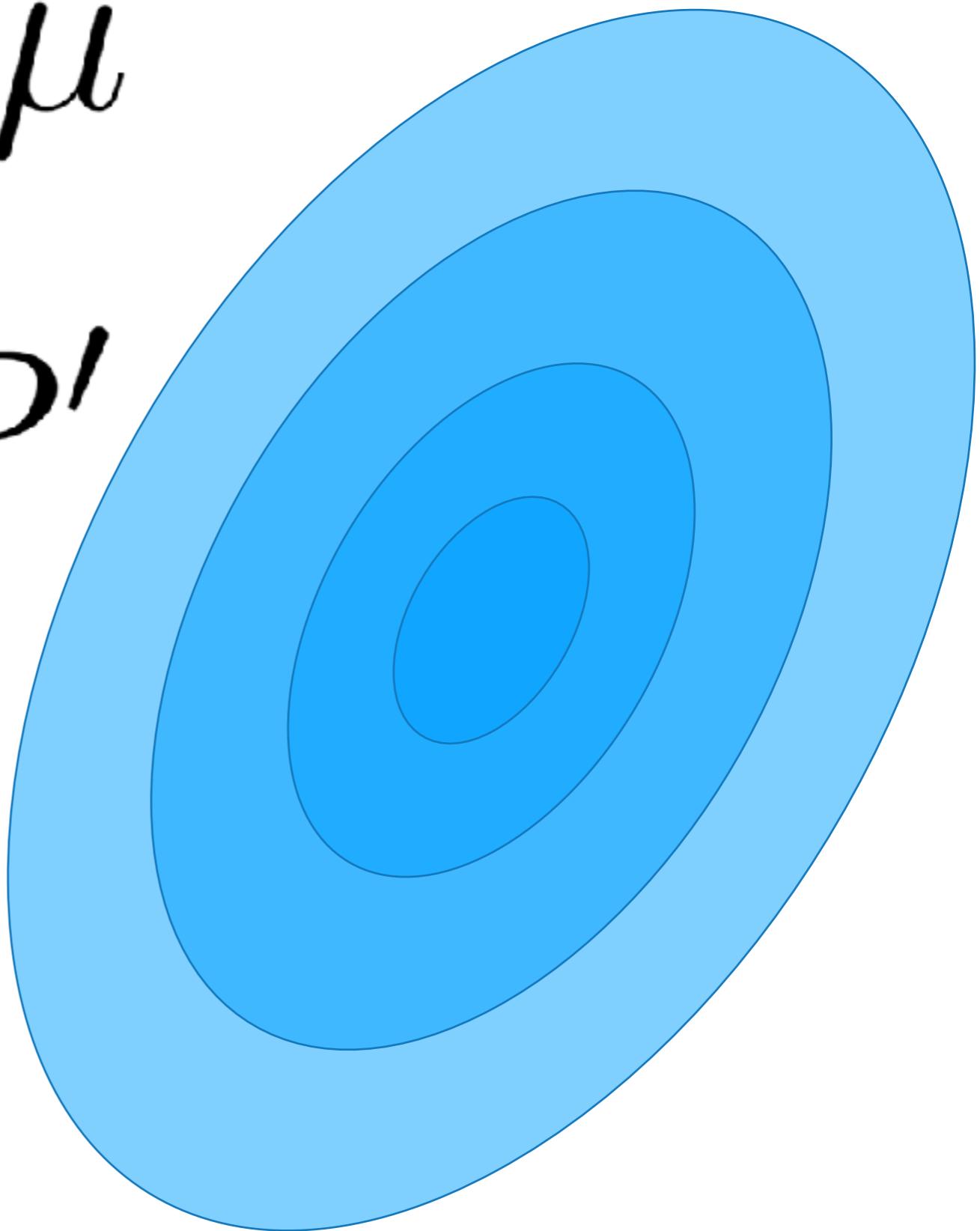


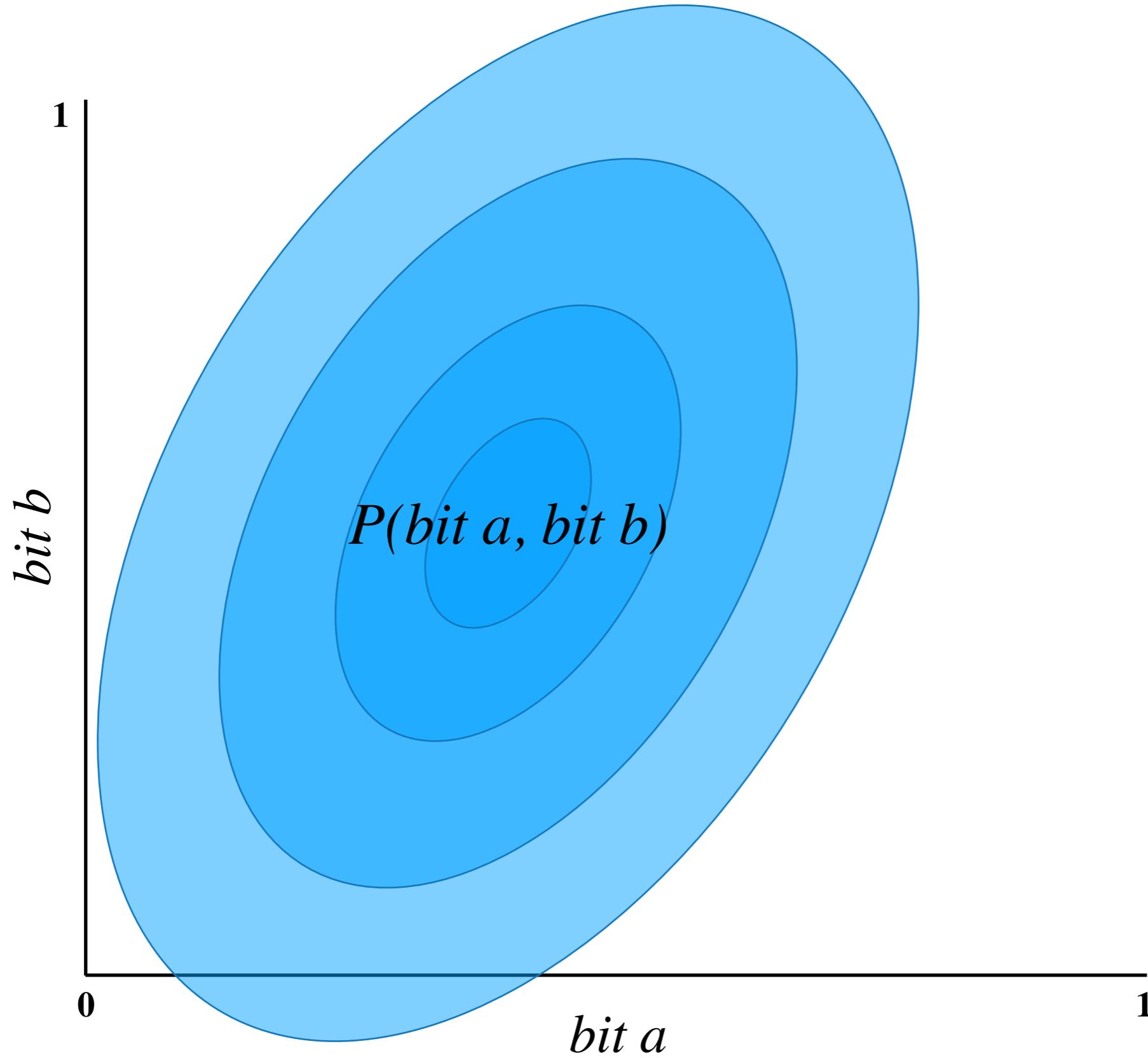
1 1

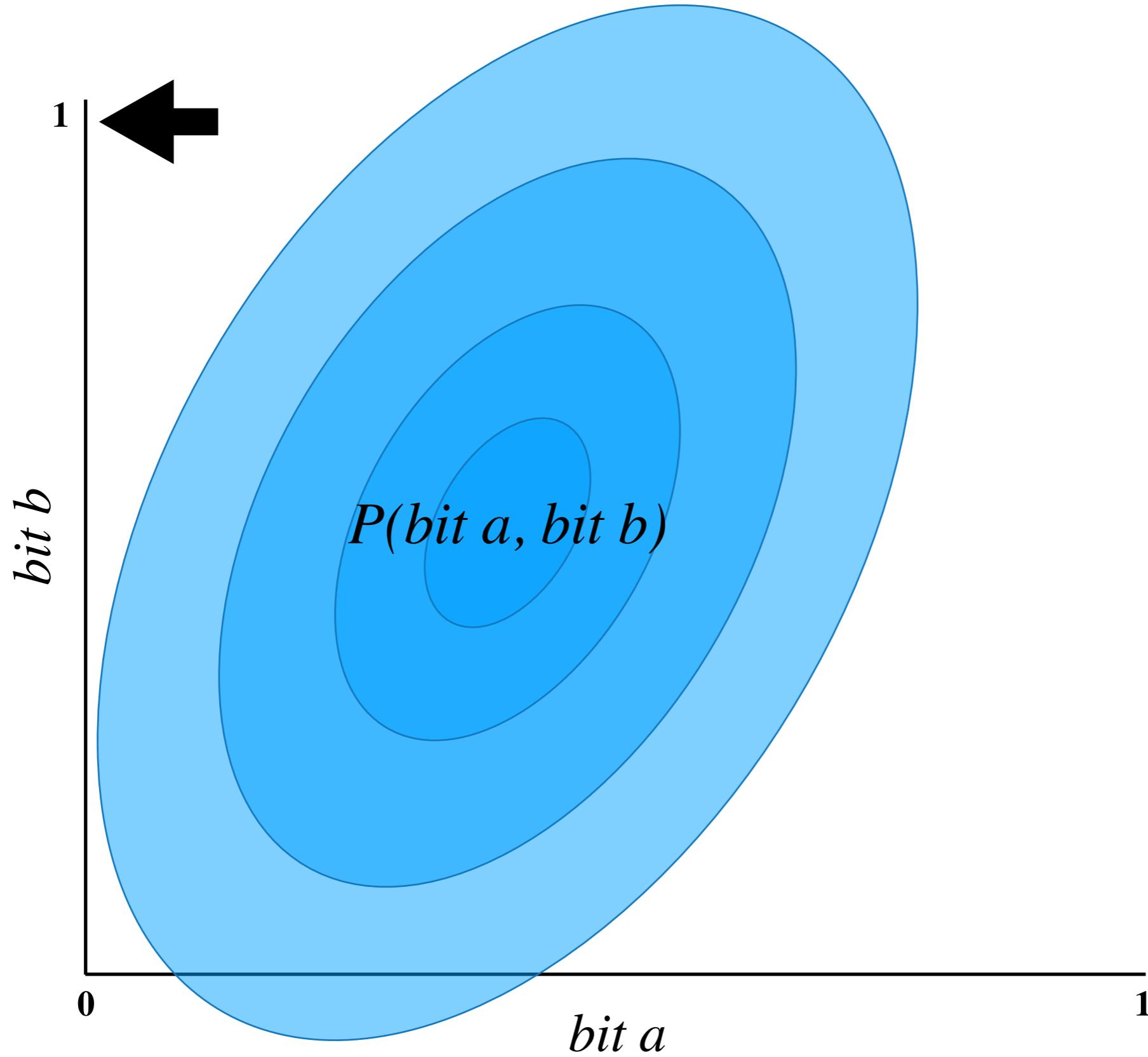
? 1 1 ? ? ? ? ? ? ?

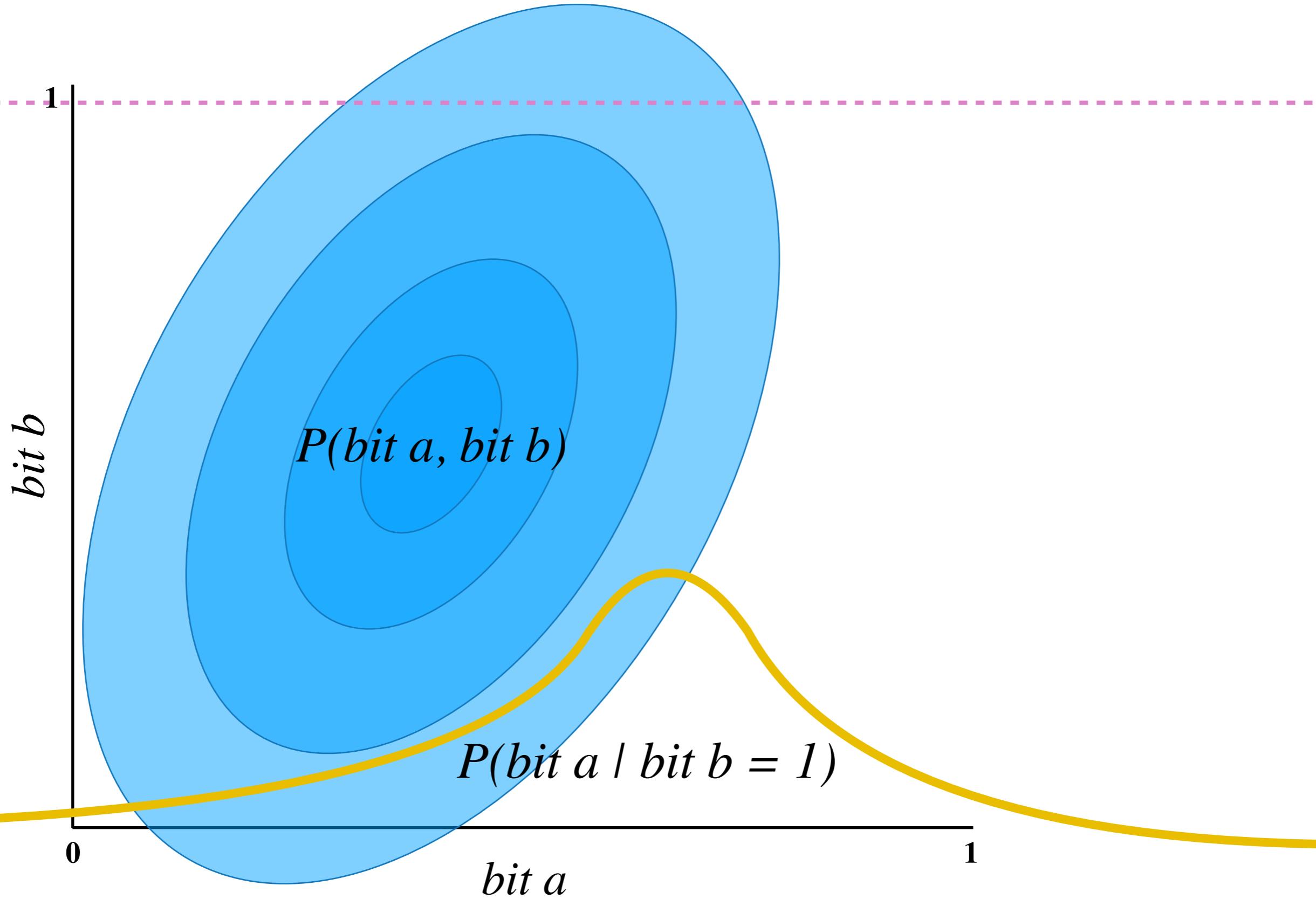
$$P' = P - \mu$$

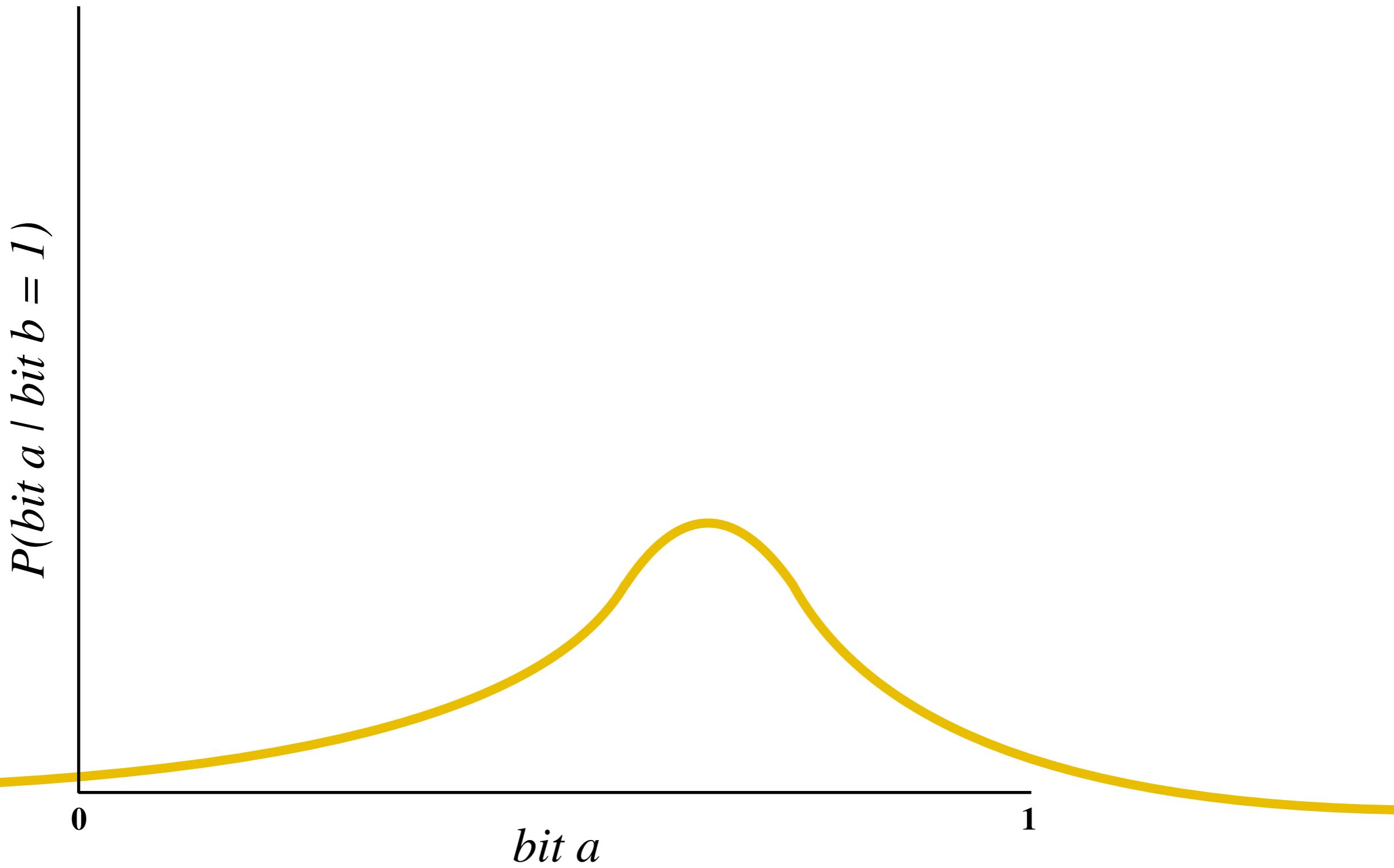
$$\Sigma = P'^T P'$$

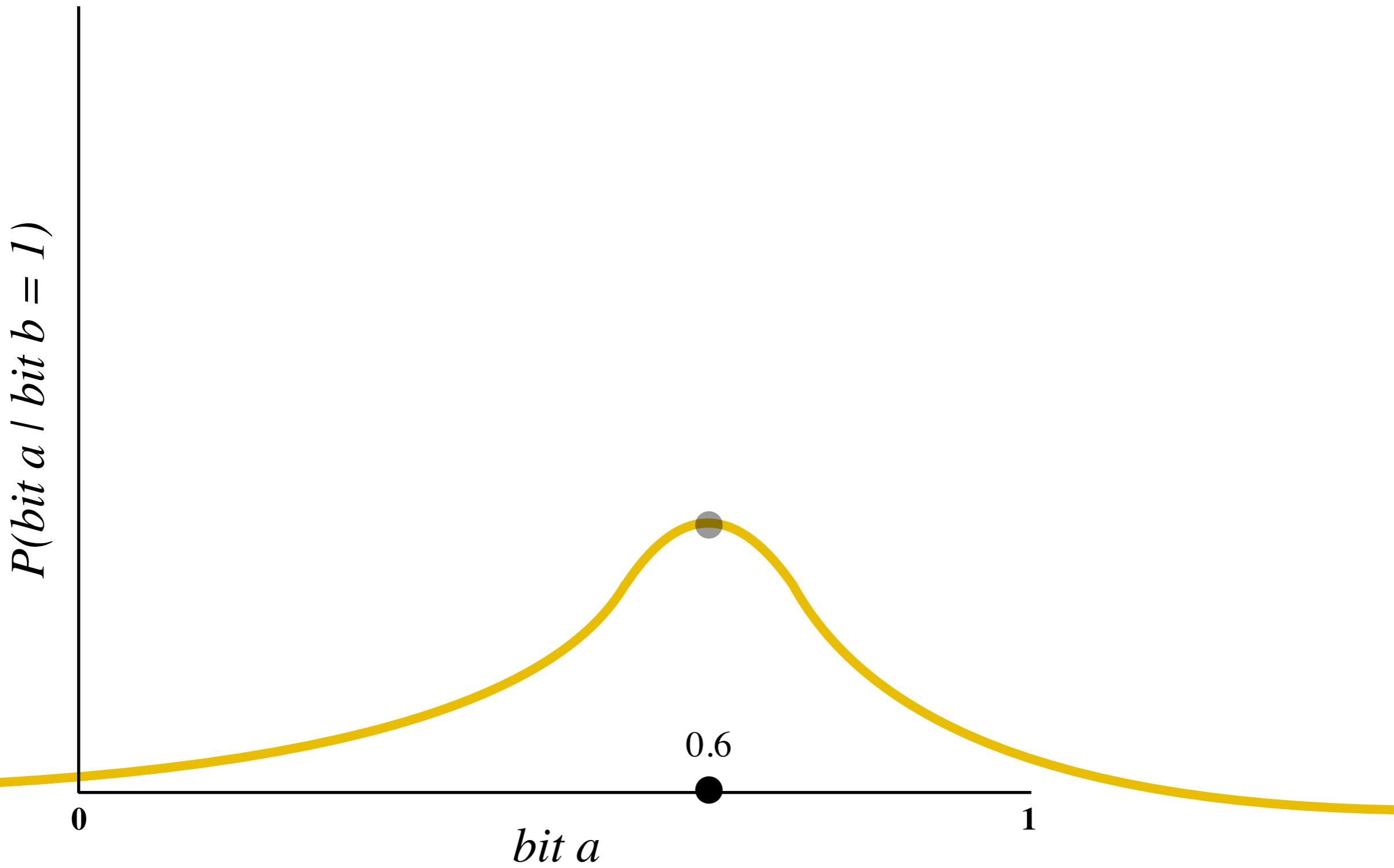










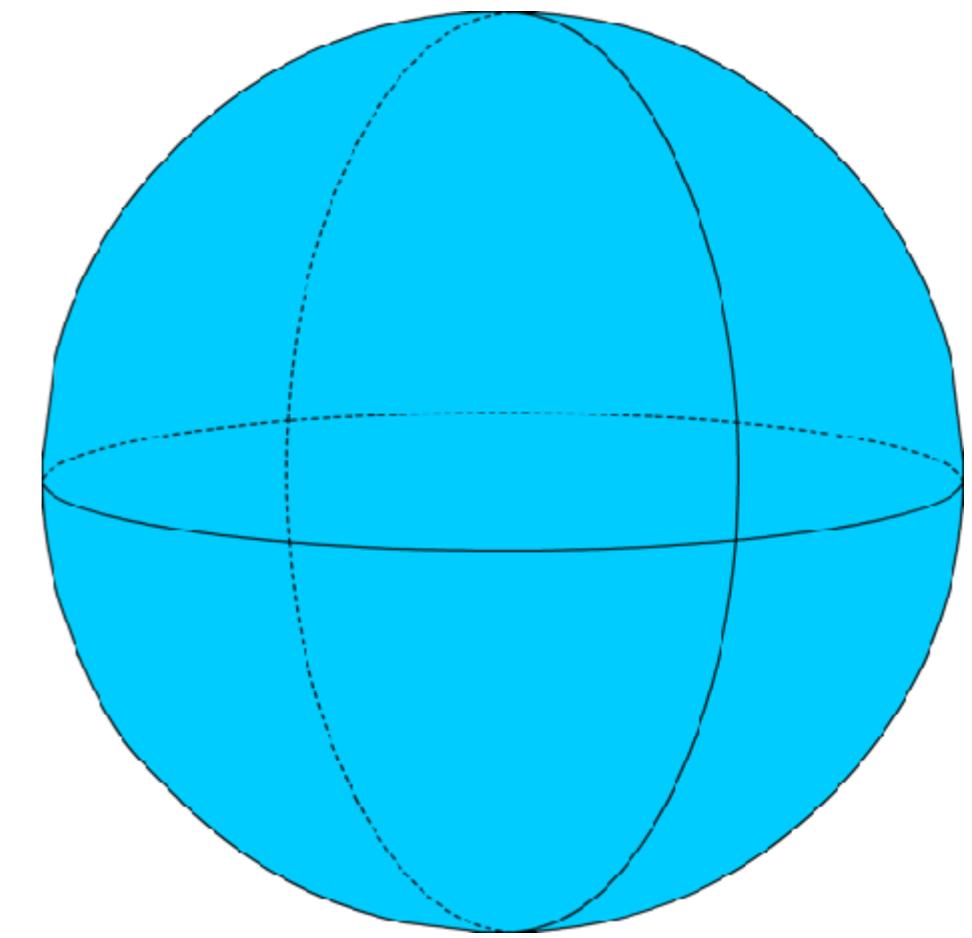




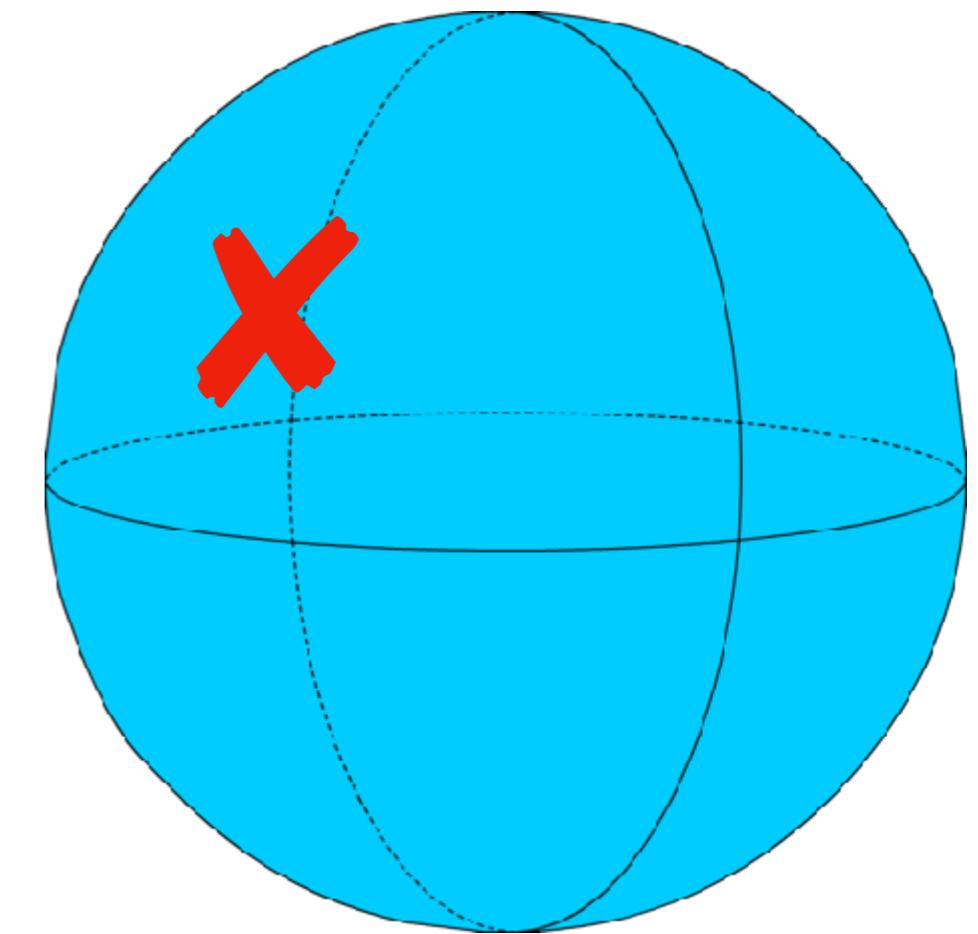
1 1

0	1	1	1	0	1	1	1	0	1	1
---	---	---	---	---	---	---	---	---	---	---

1
1
⋮



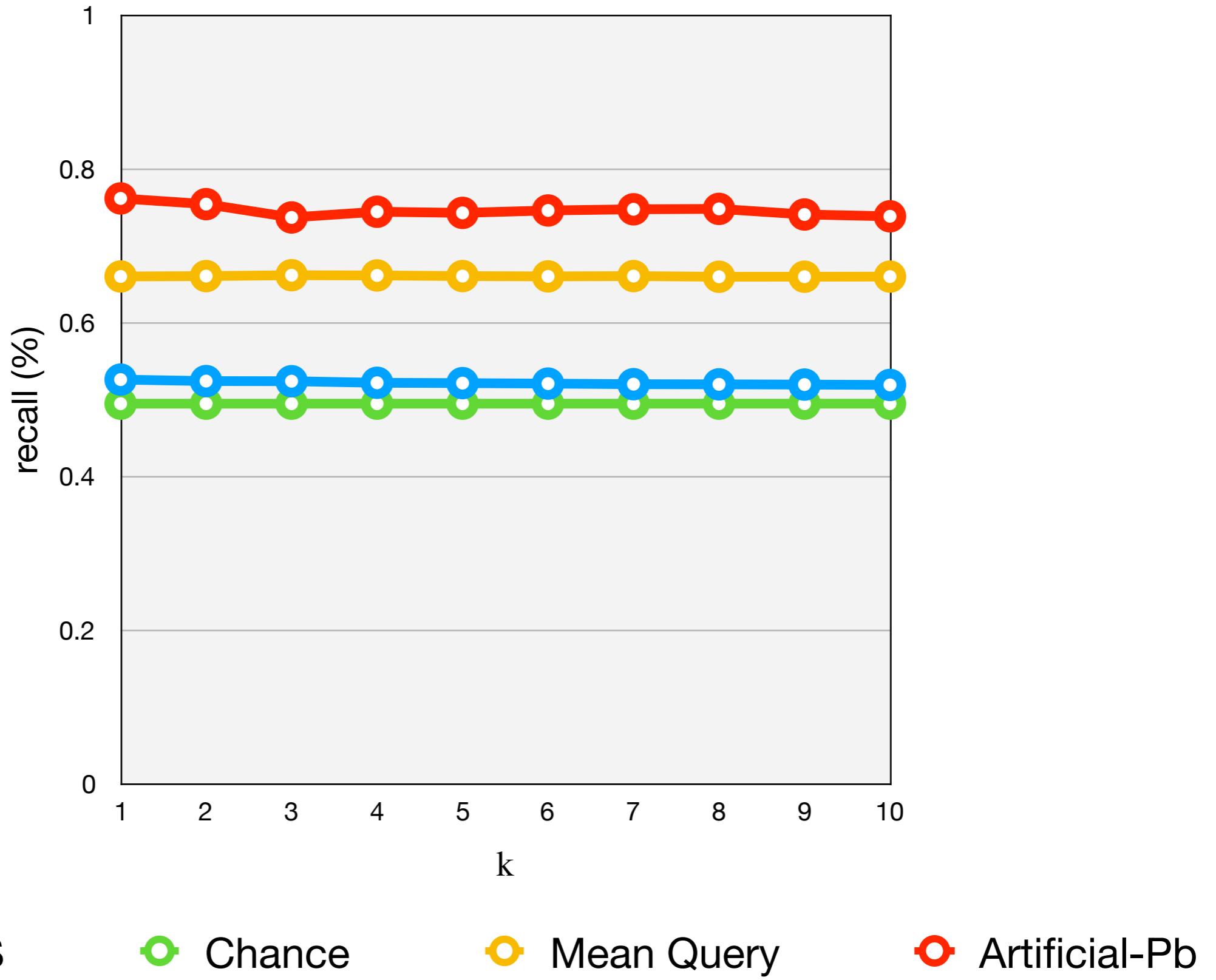
0
1
1
:
0
1
1
0



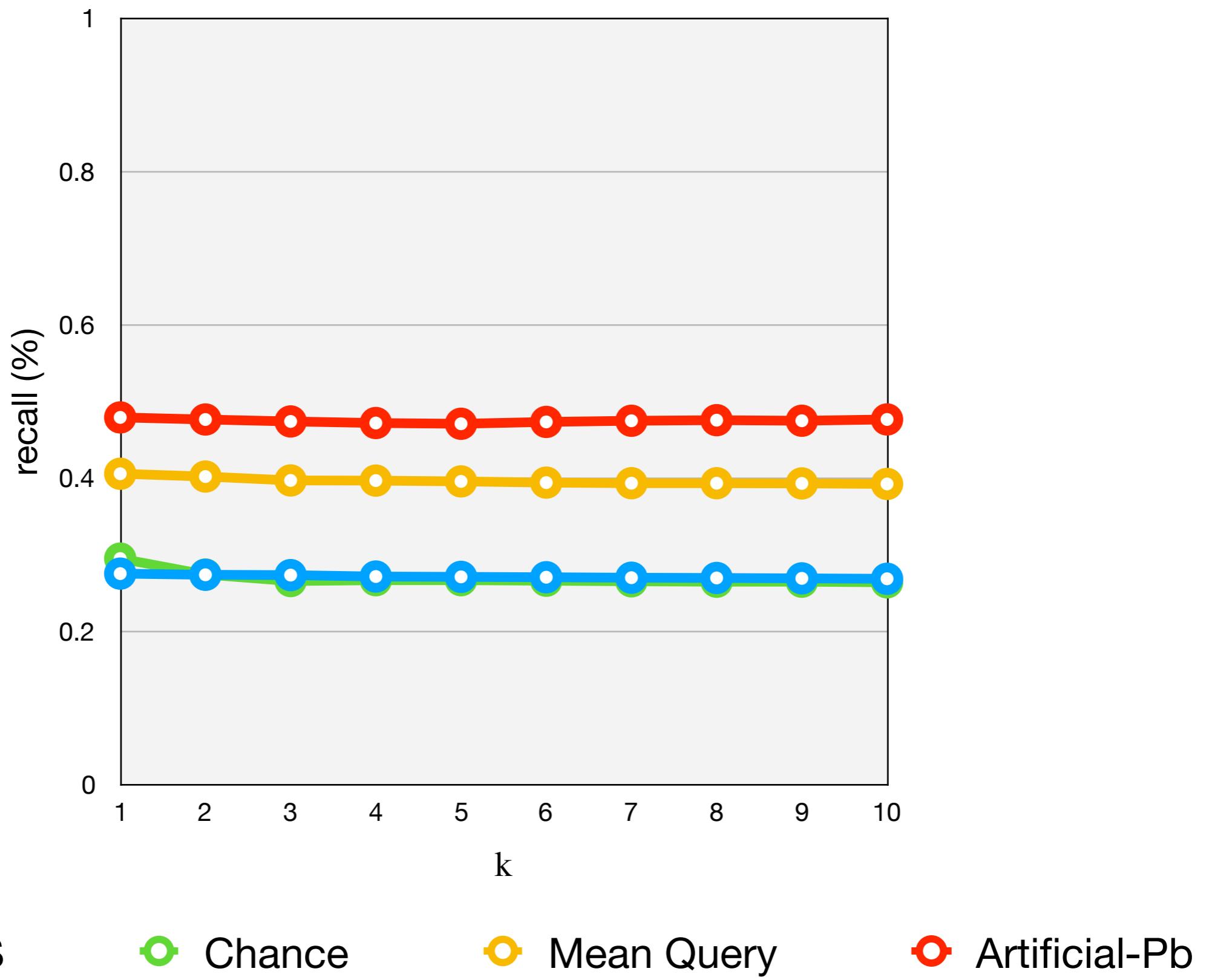
$$recall@k = \frac{1}{N} \sum_{i=1}^N recall_k(R_{1..k})$$

$$recall_k(X) = \frac{|X \cap X_{tp}|}{|X|}$$

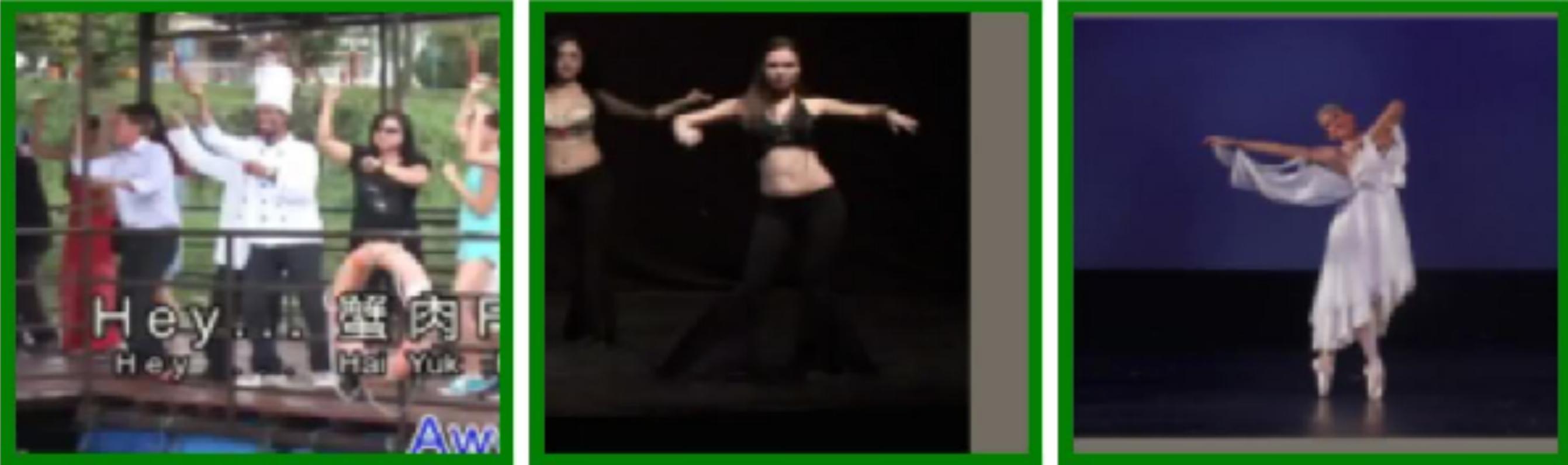
recall@k (1st order)

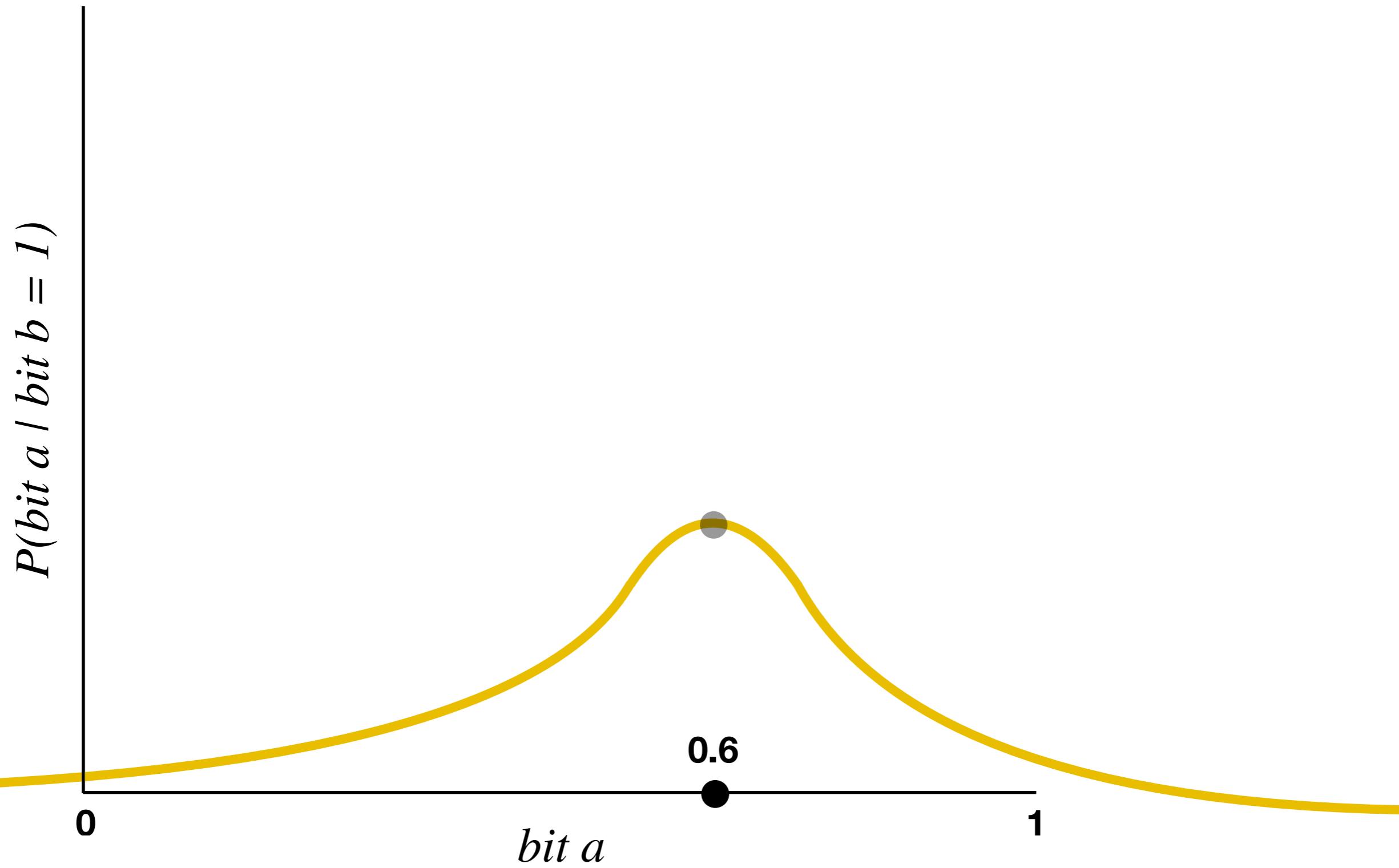


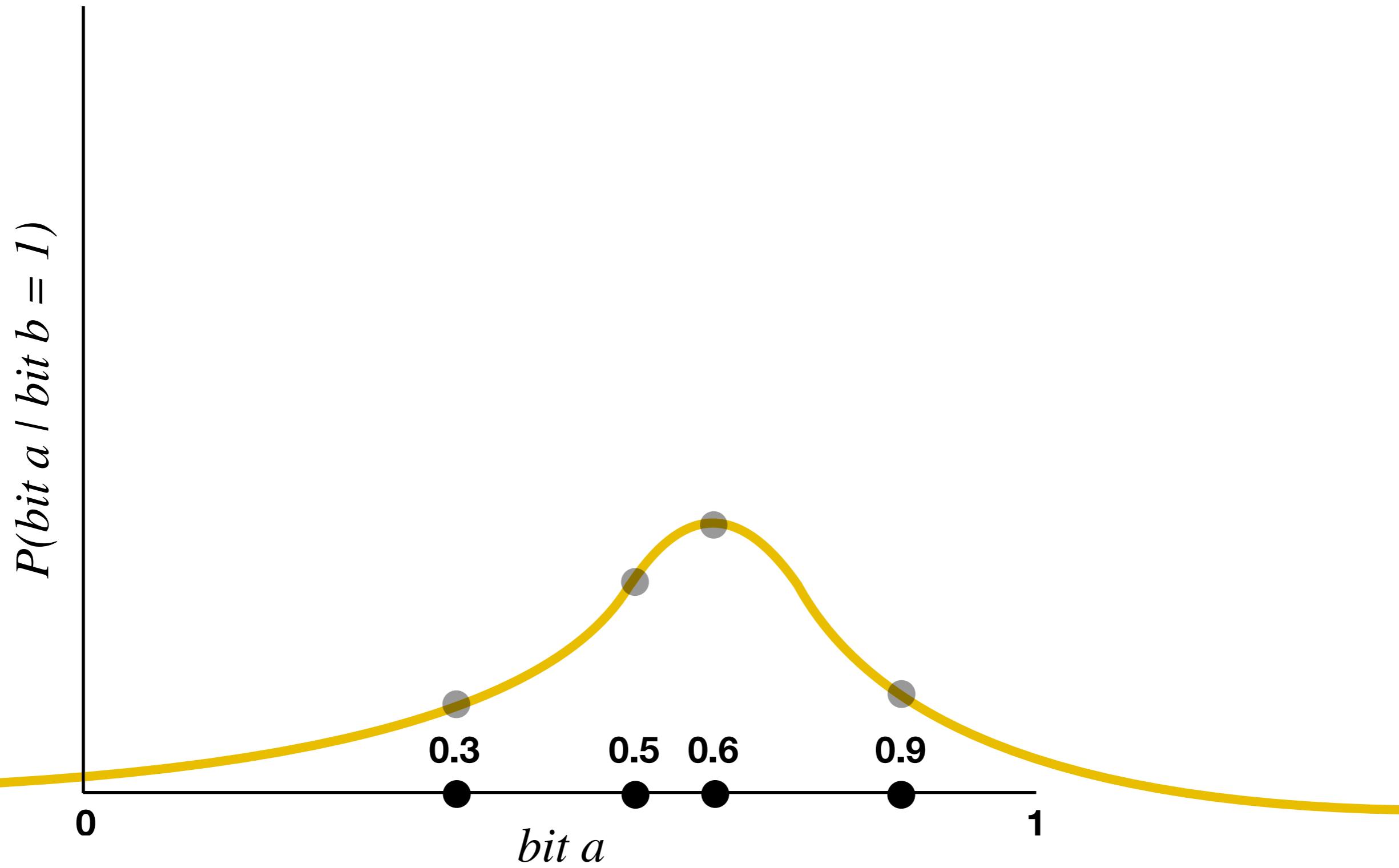
recall@k (2nd order)

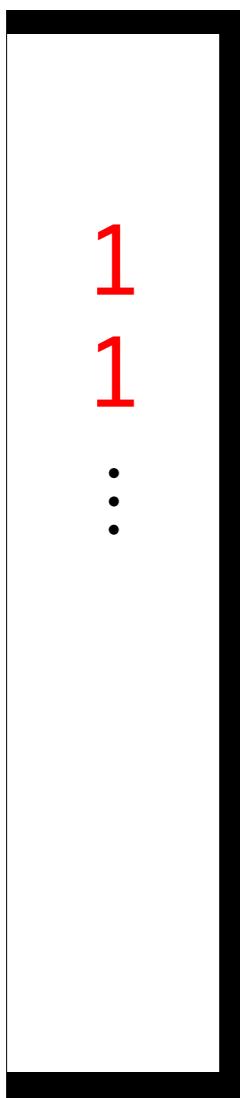


right shoulder is bent



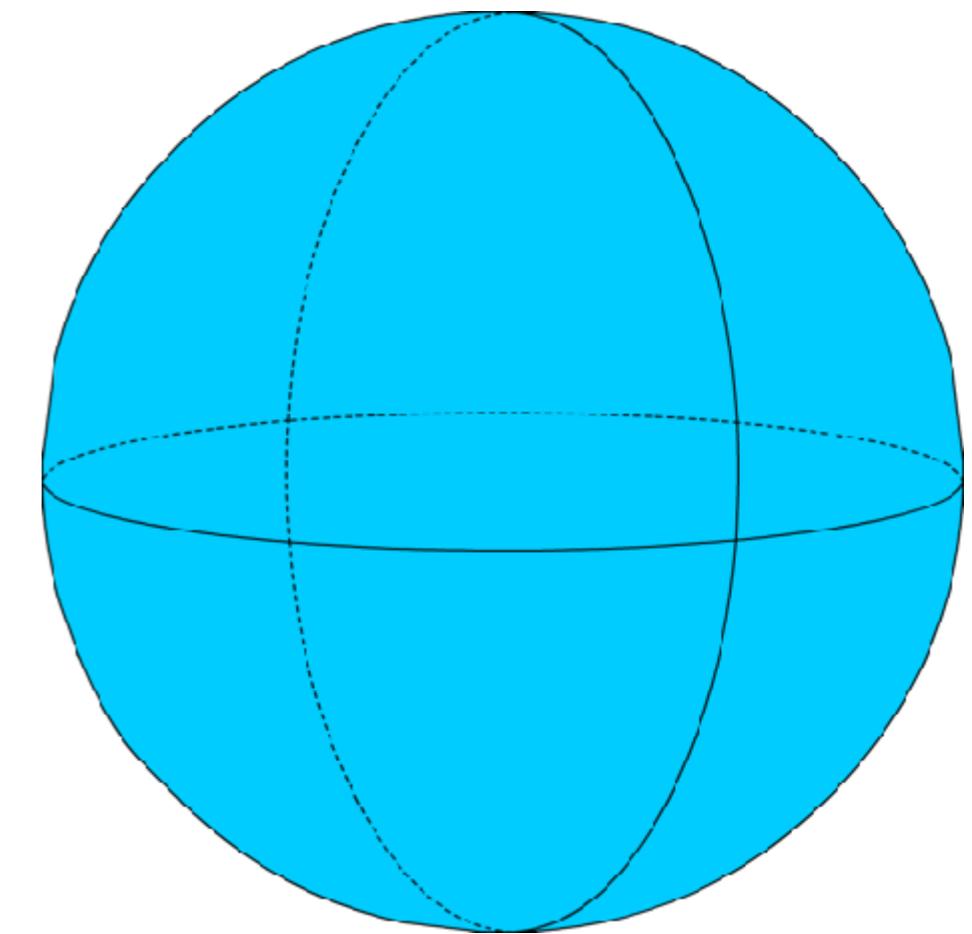




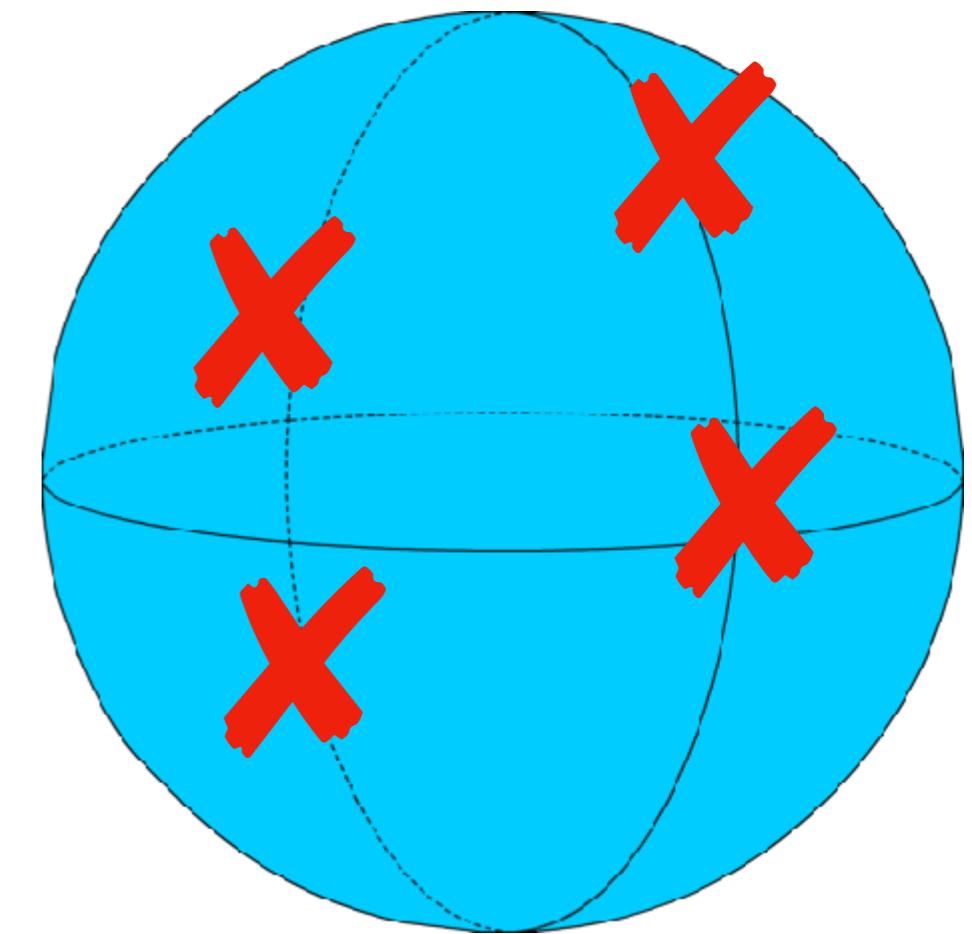


1
1

⋮



1
1
1
:
1
0
1
1



right shoulder is bent



Mahalanobis Distance

mean recall@1

0.0

76.2%

0.1

73.1%

1.0

71.1%

10.0

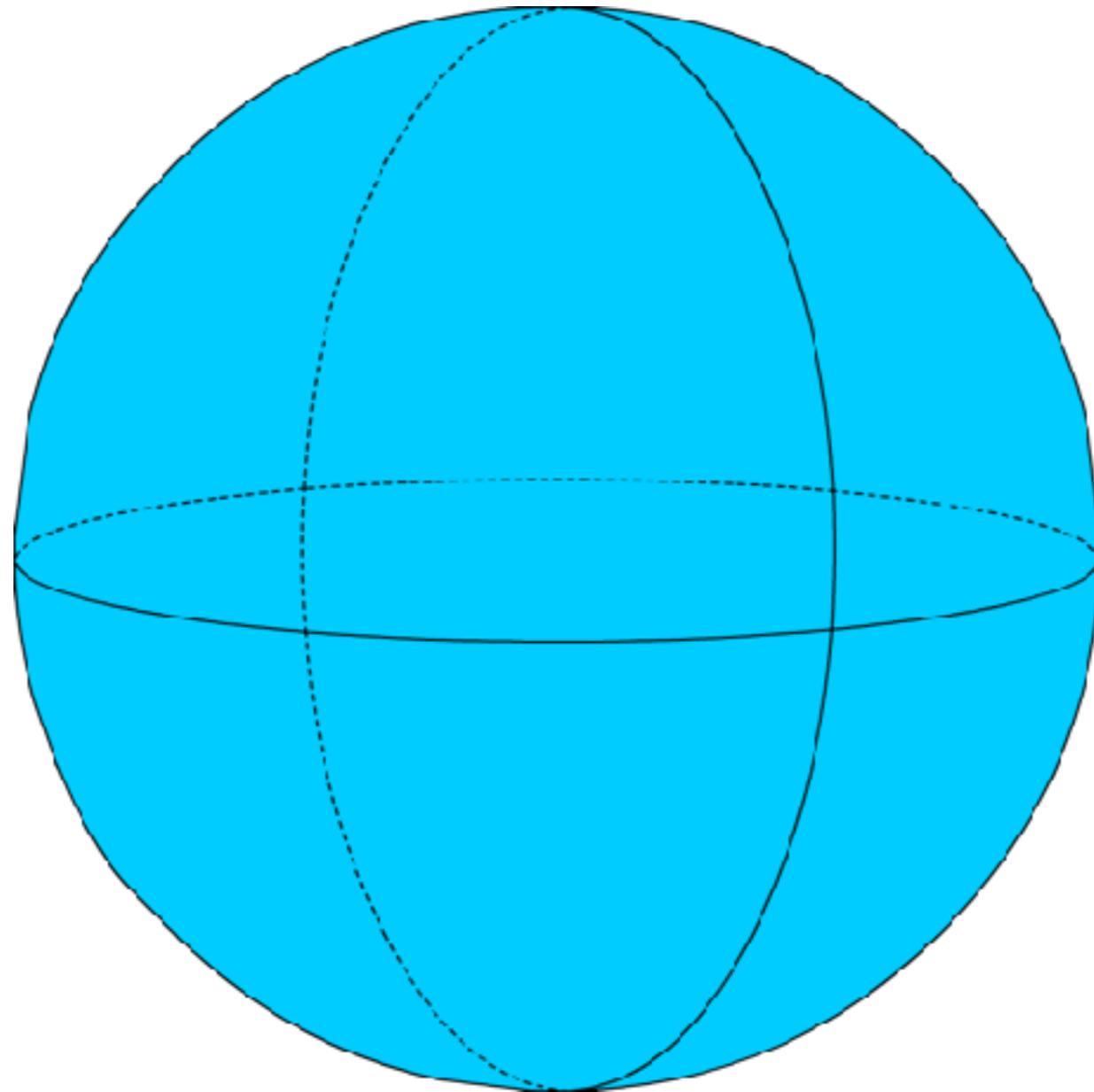
64.7%

Baseline

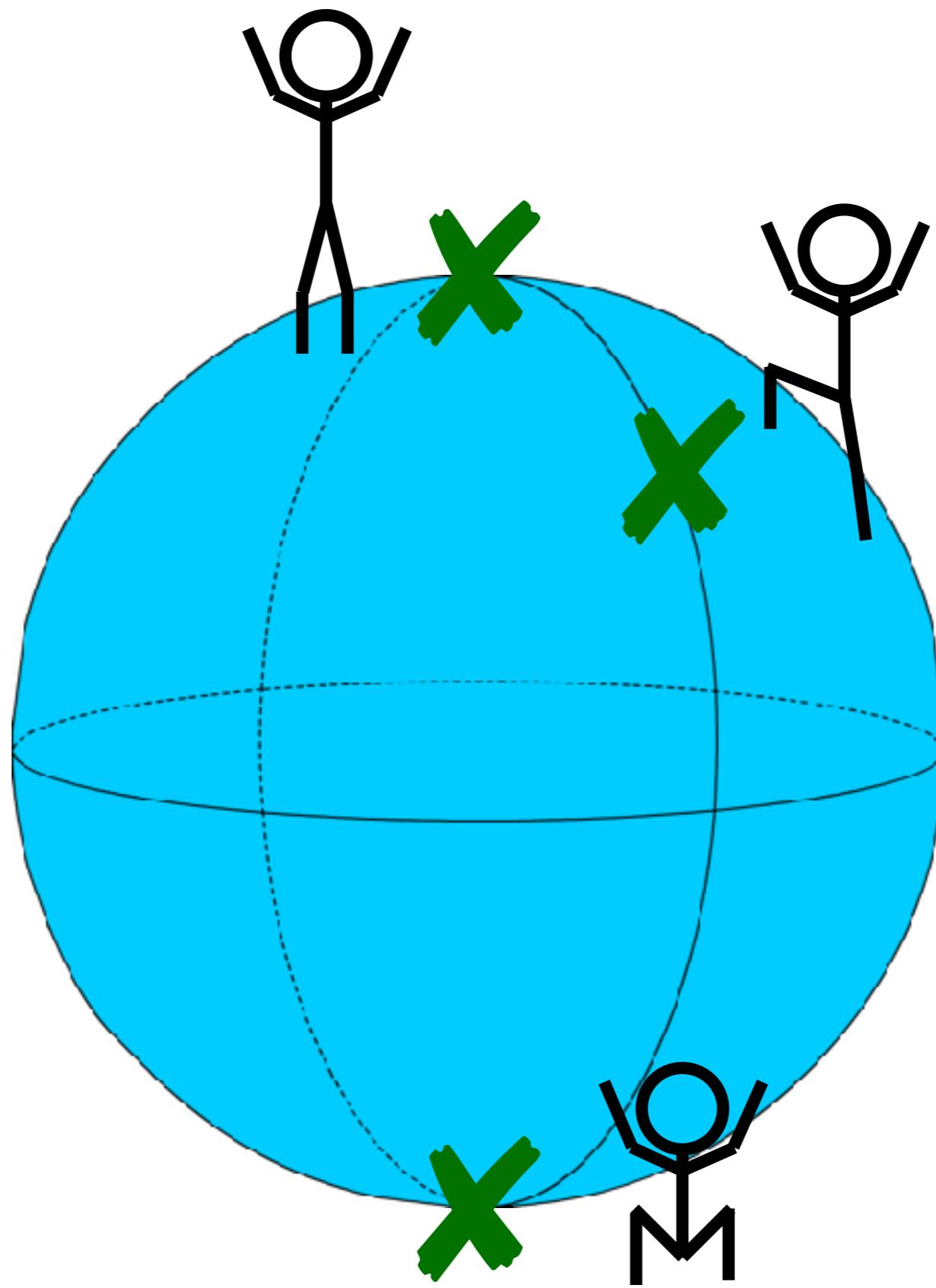
Mean Query

66.1%

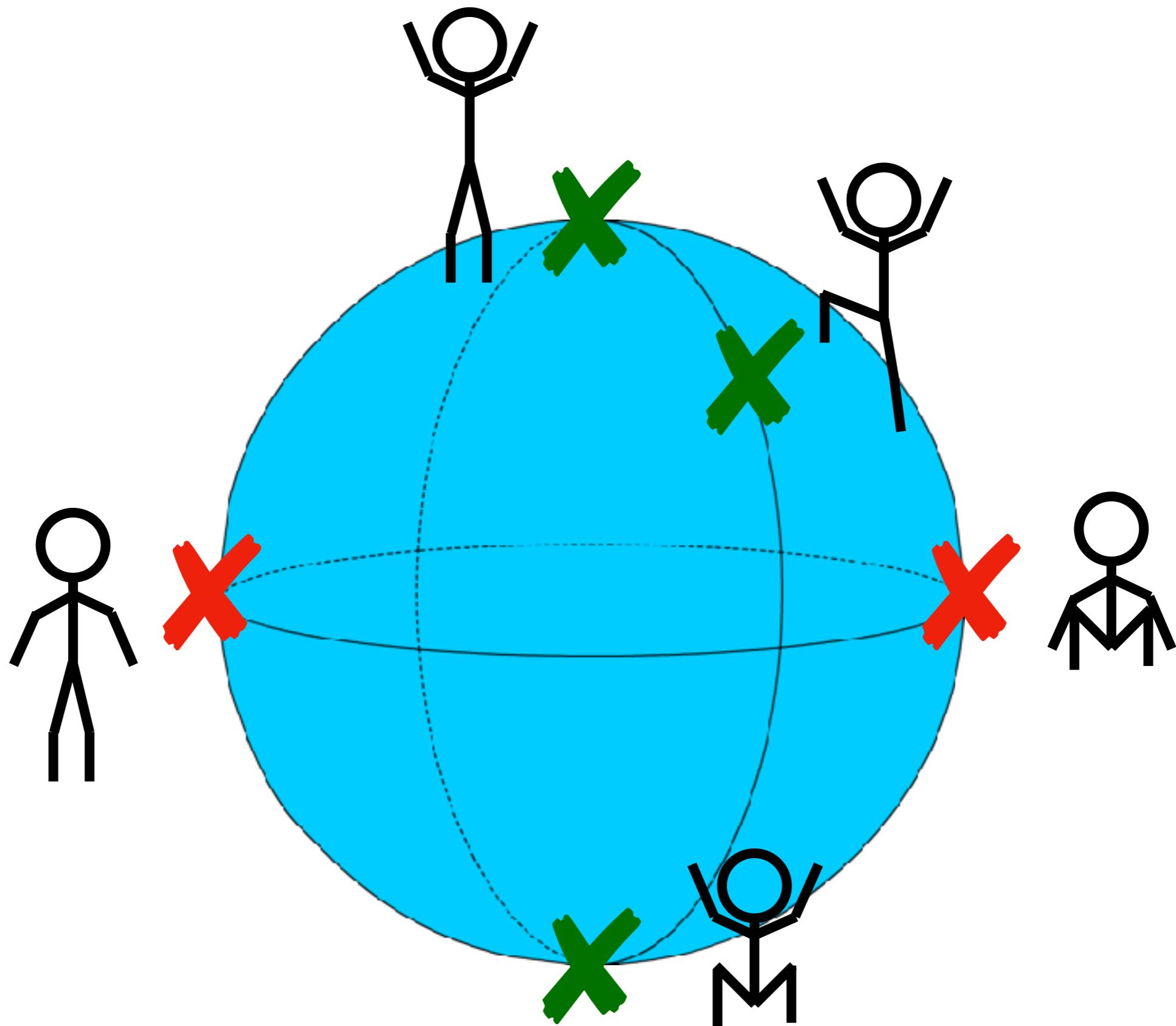
“Left hand above head; right hand above head.”



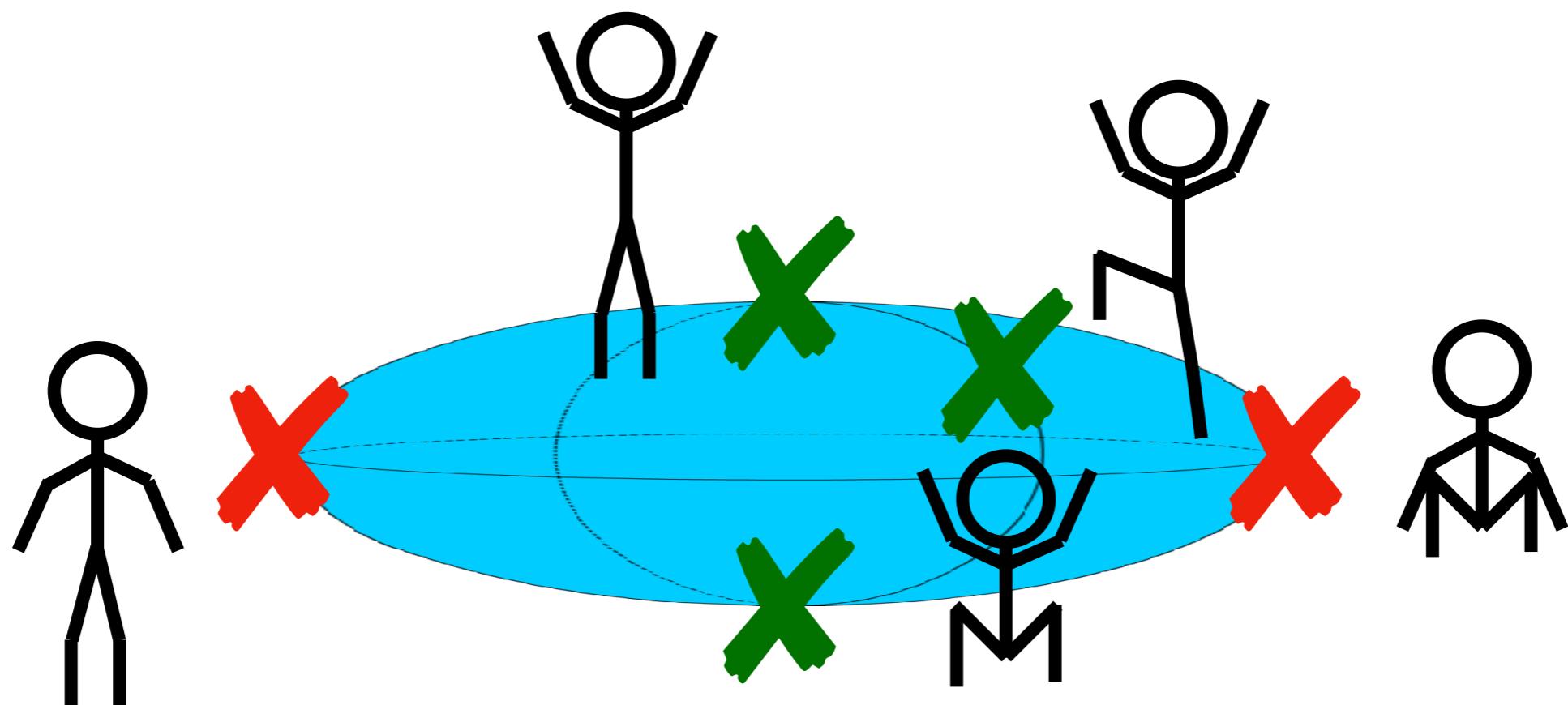
“Left hand above head; right hand above head.”



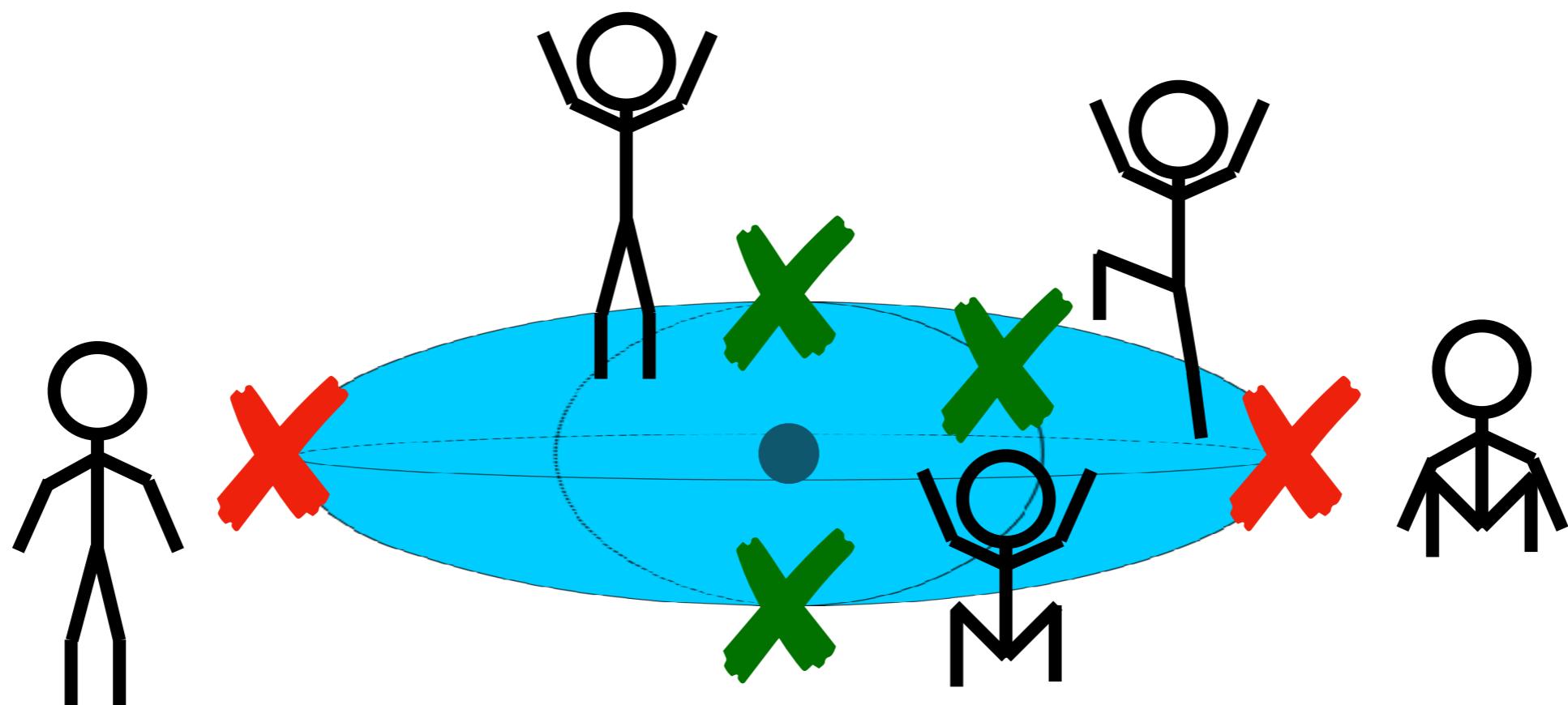
“Left hand above head; right hand above head.”



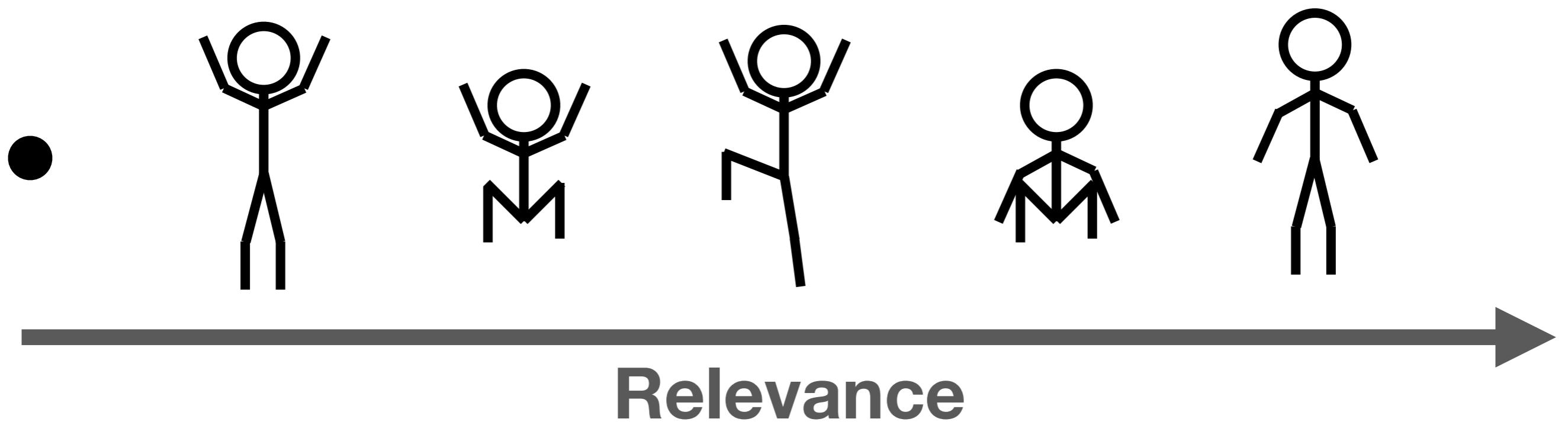
“Left hand above head; right hand above head.”

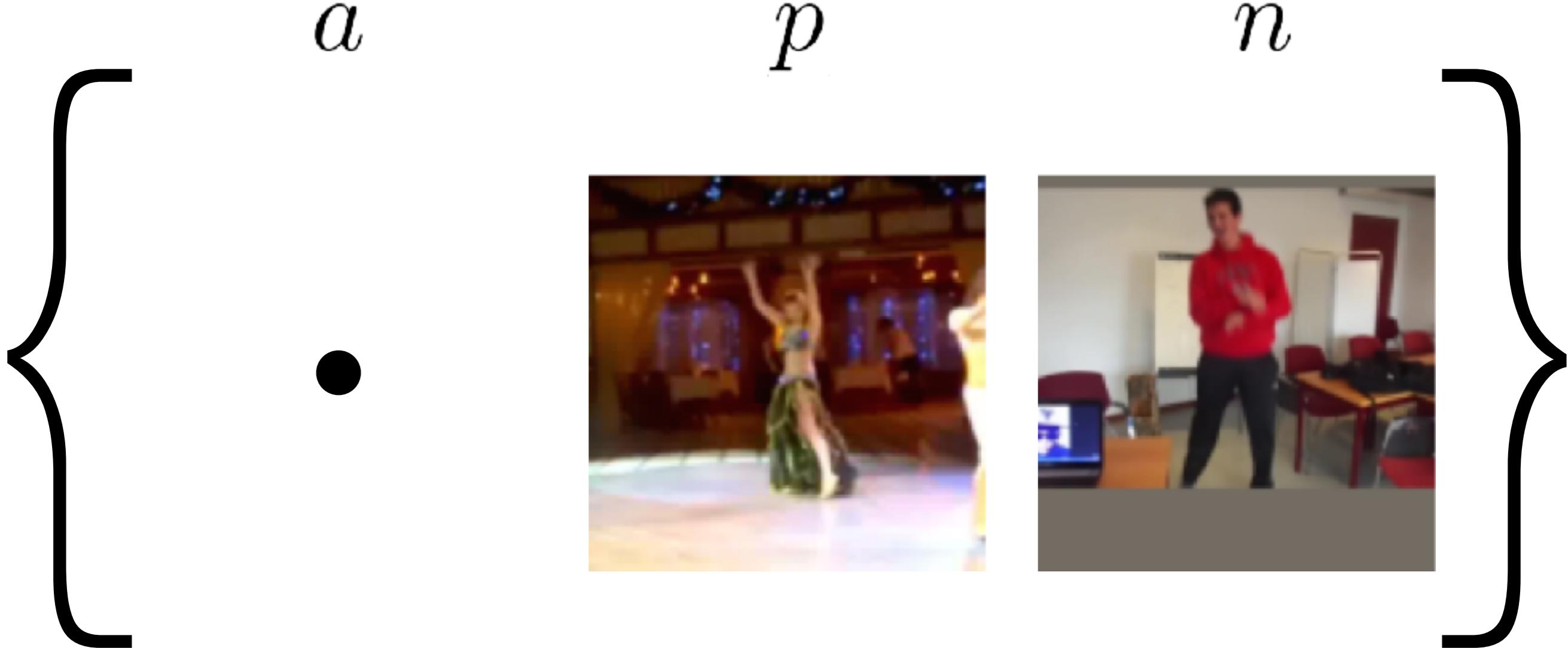


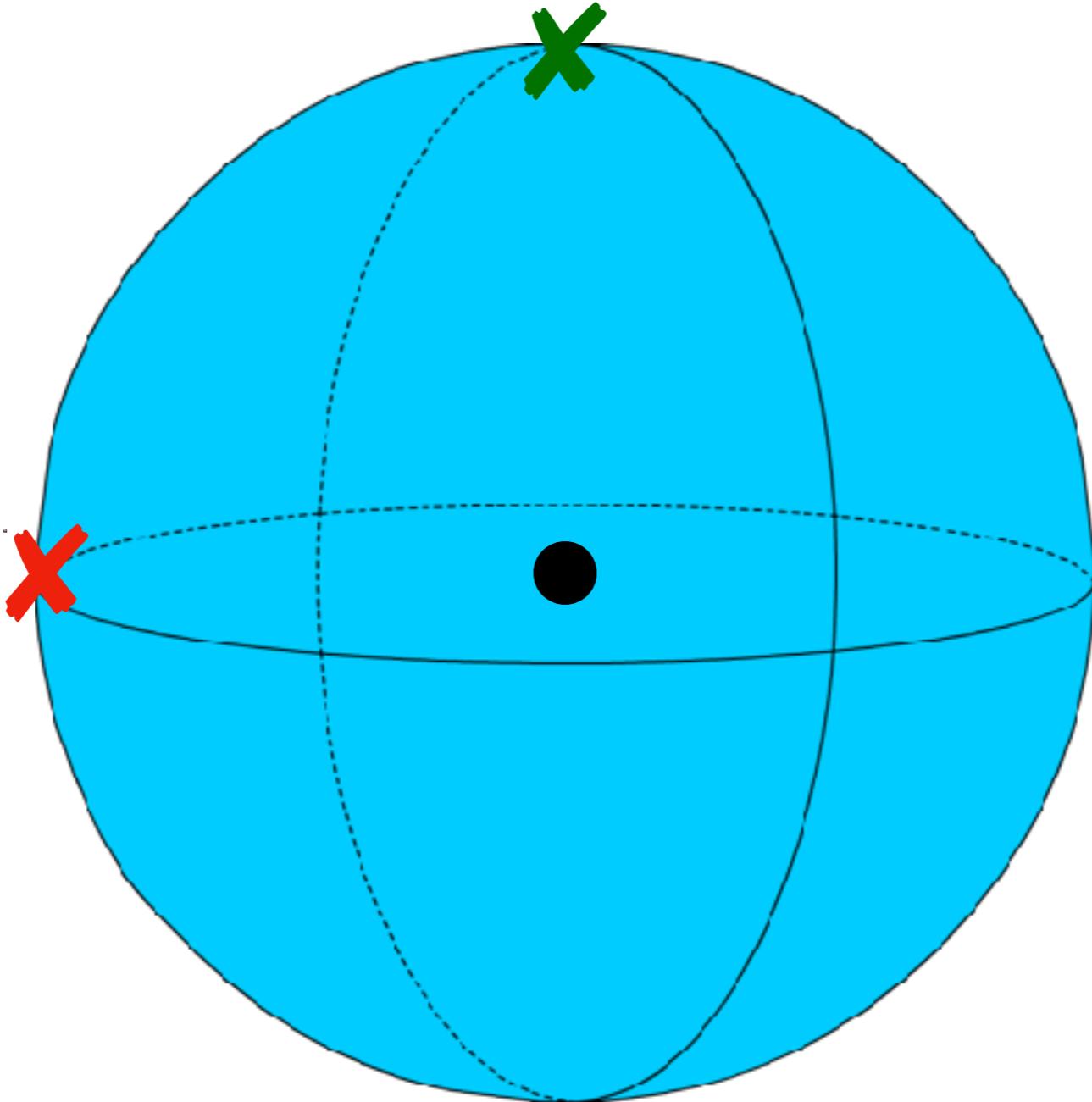
“Left hand above head; right hand above head.”



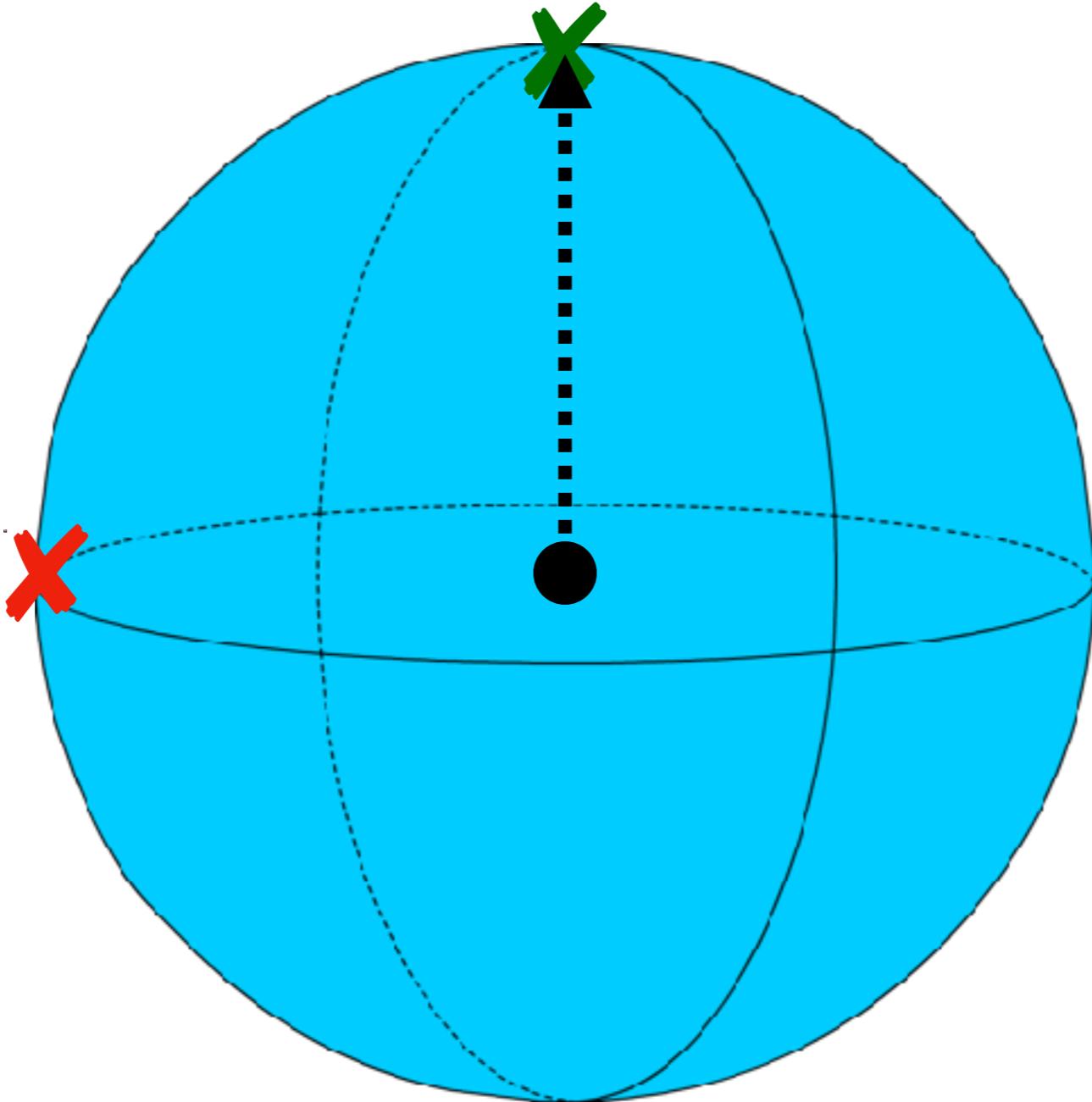
“Left hand above head; right hand above head.”



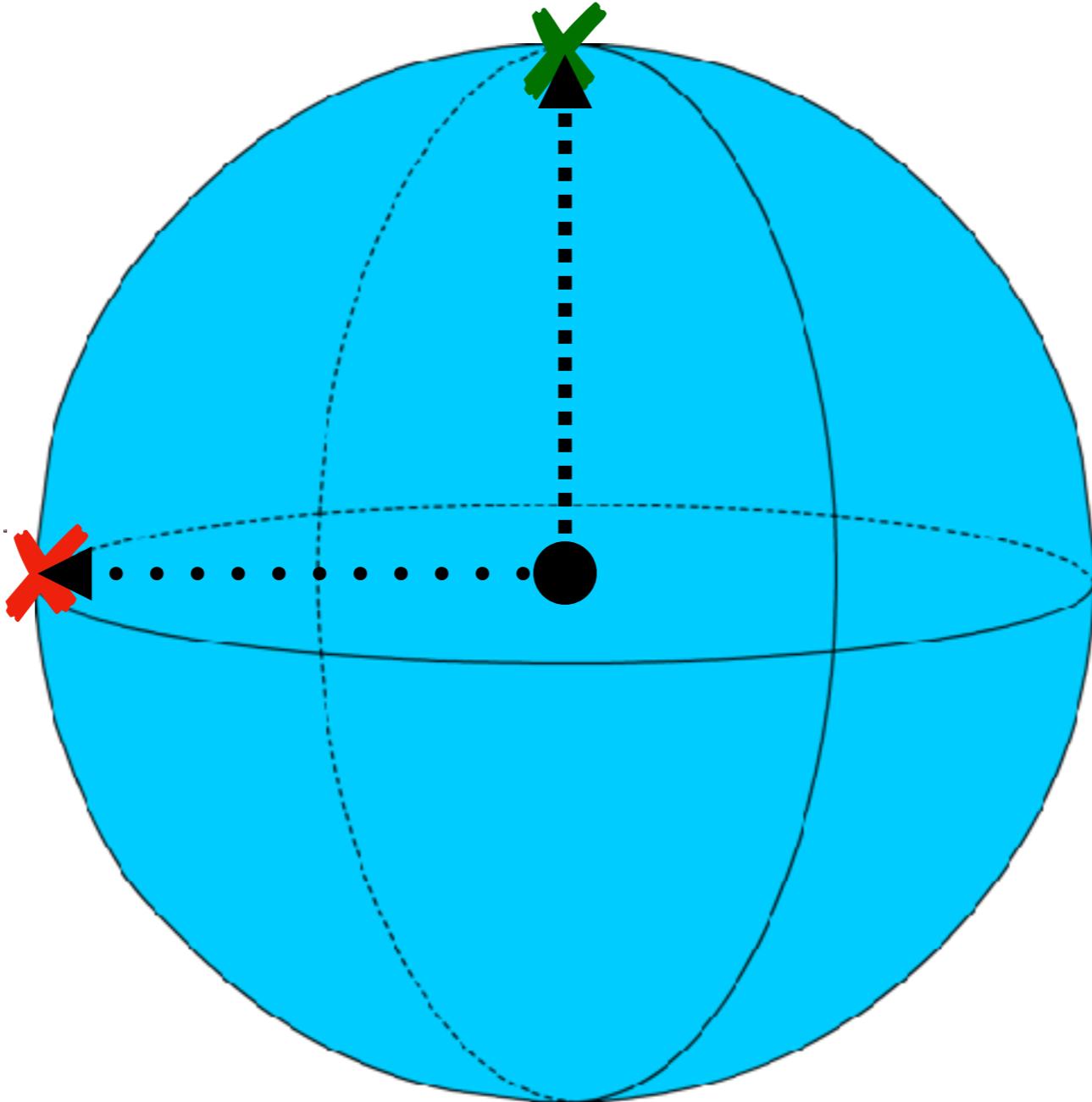




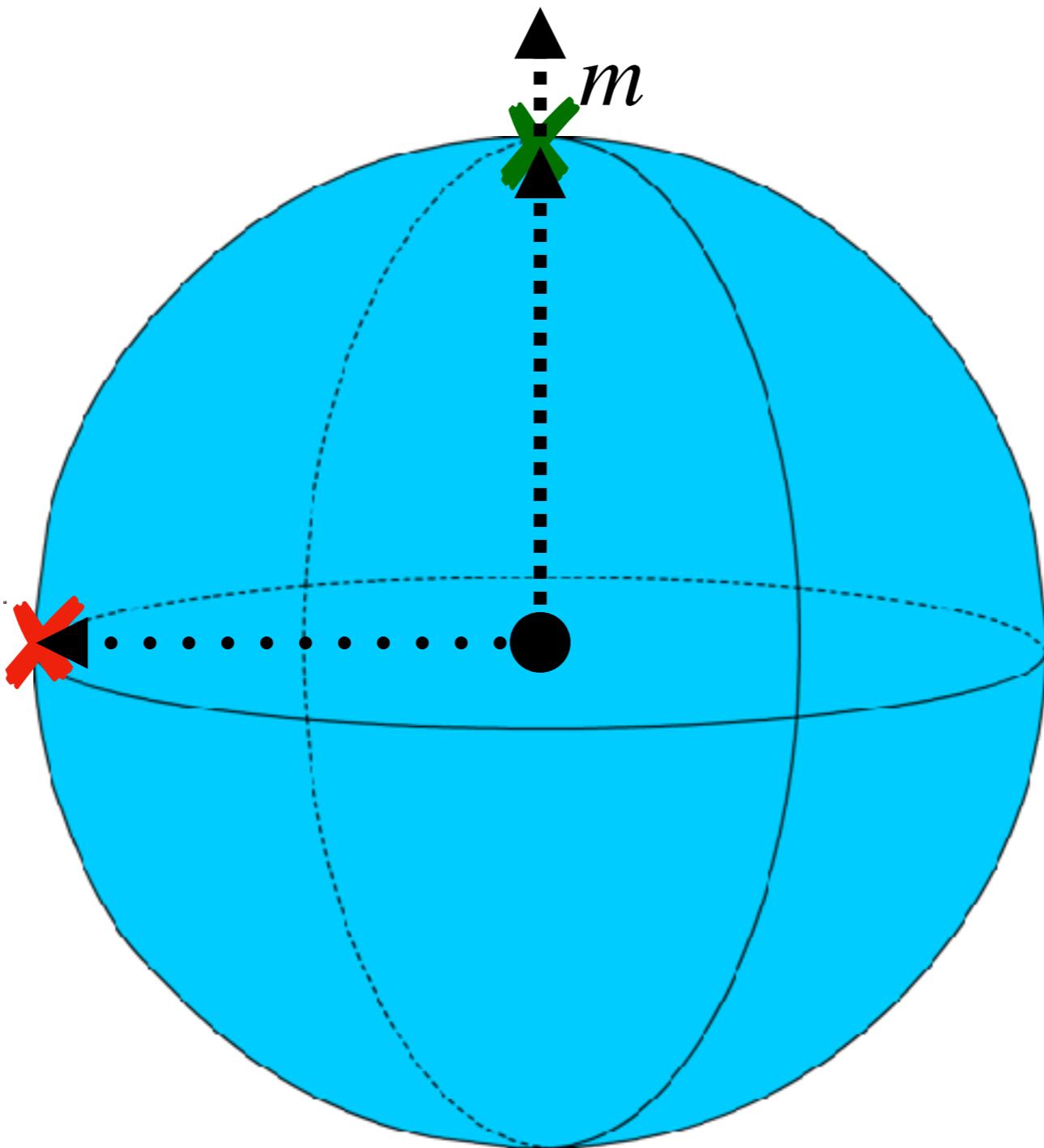
$$\mathcal{L}_{triplet} = [||f(a) - f(p)||_2^2 - ||f(a) - f(n)||_2^2 + m]_+$$



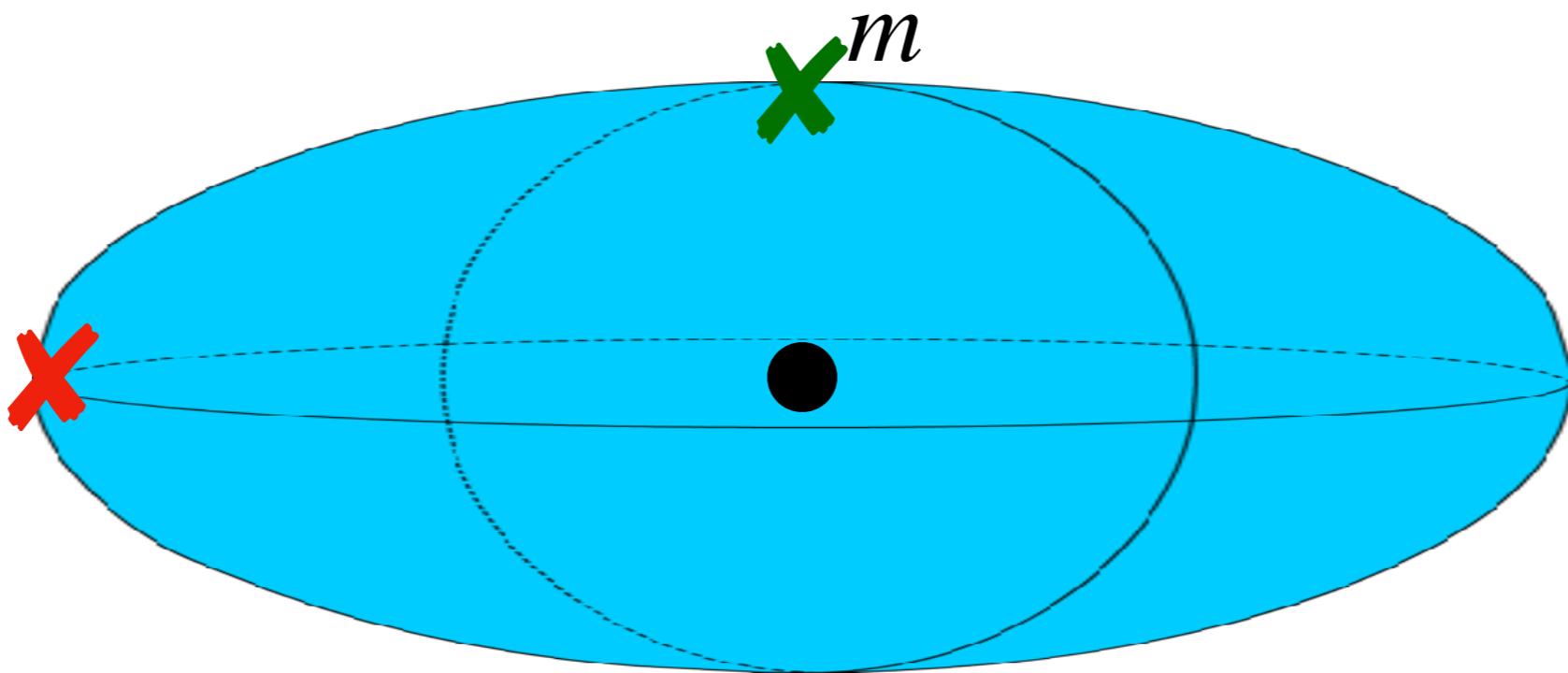
$$\mathcal{L}_{triplet} = [||f(a) - f(p)||_2^2 - ||f(a) - f(n)||_2^2 + m]_+$$



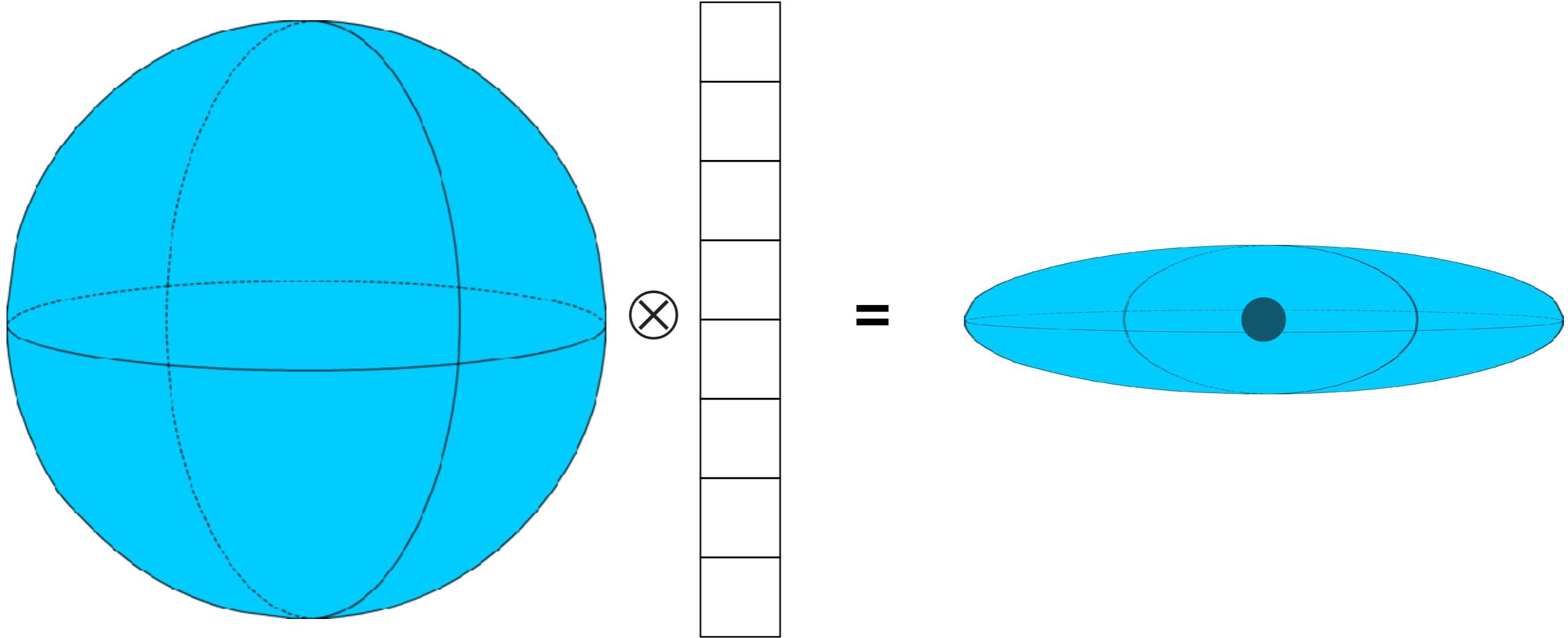
$$\mathcal{L}_{triplet} = [||f(a) - f(p)||_2^2 - ||f(a) - f(n)||_2^2 + m]_+$$



$$\mathcal{L}_{triplet} = [||f(a) - f(p)||_2^2 - ||f(a) - f(n)||_2^2 + m]_+$$



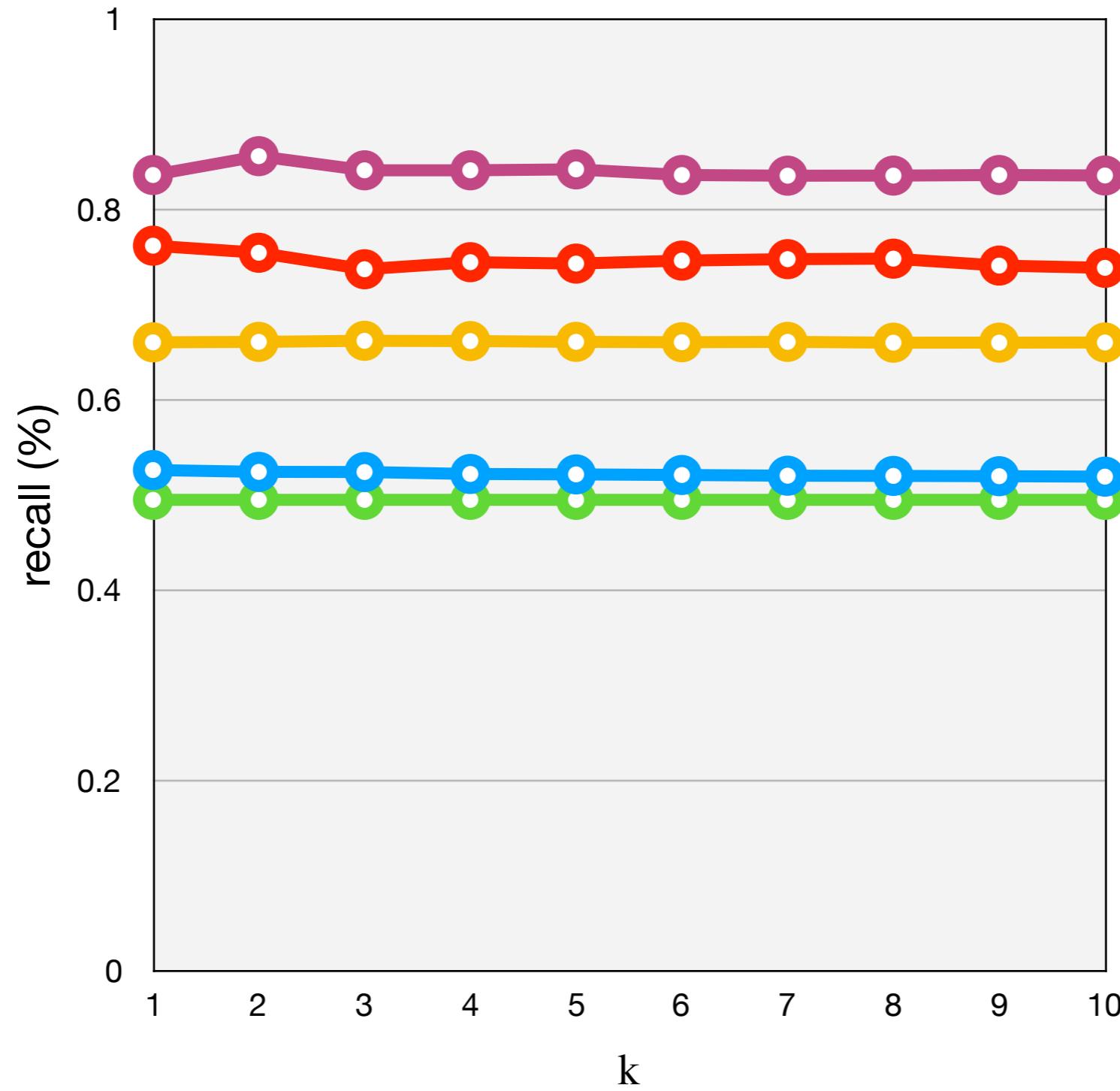
$$\mathcal{L}_{triplet} = [||f(a) - f(p)||_2^2 - ||f(a) - f(n)||_2^2 + m]_+$$



“right wrist **is above** head”

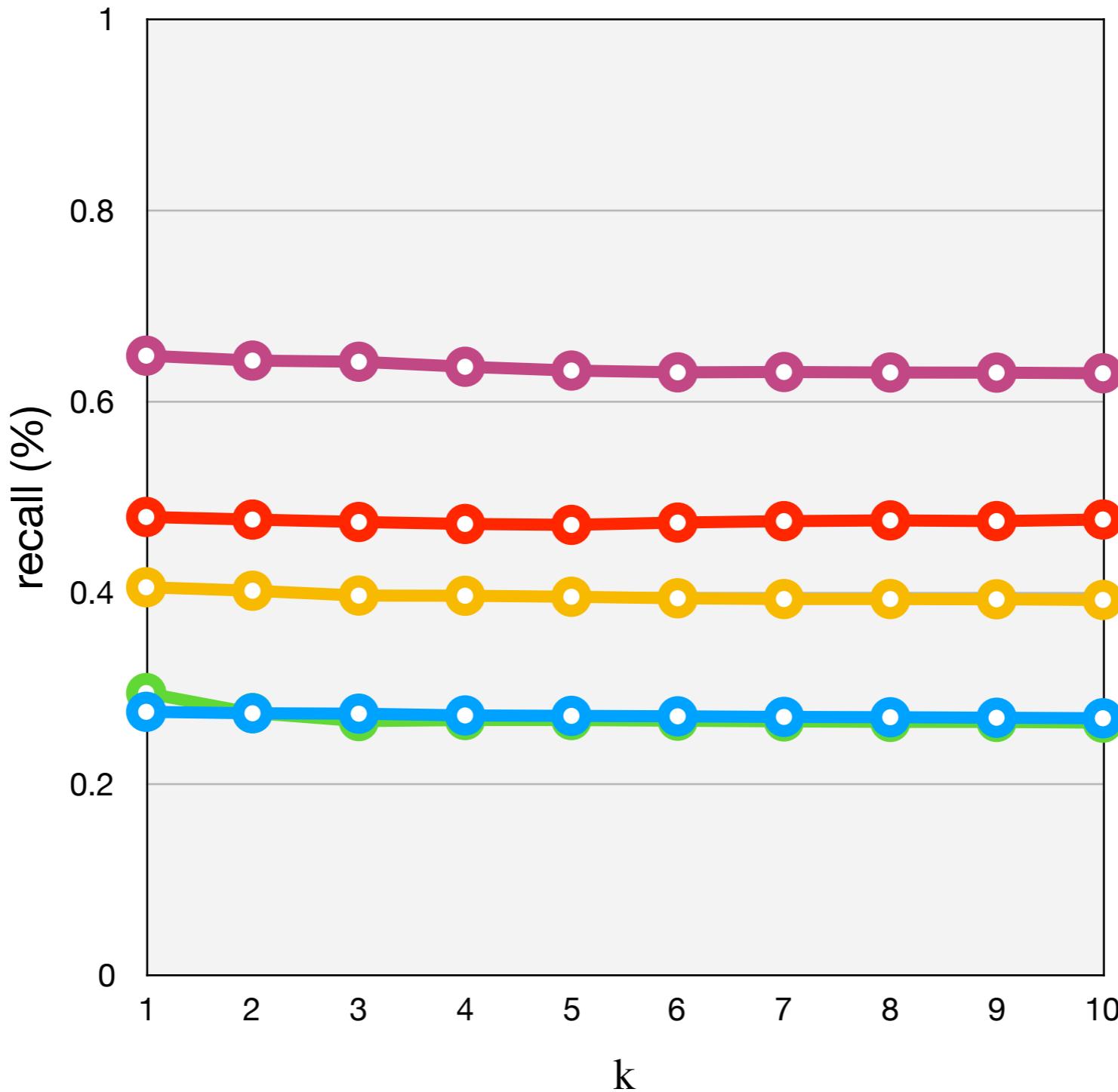
$$e'_{emb} = e_{emb} v_{mask}$$

recall@k (1st order)



○ VGG-S ○ Chance ○ Mean Query ○ Artificial-Pb ○ Query-Aware Mask

recall@k (2nd order)

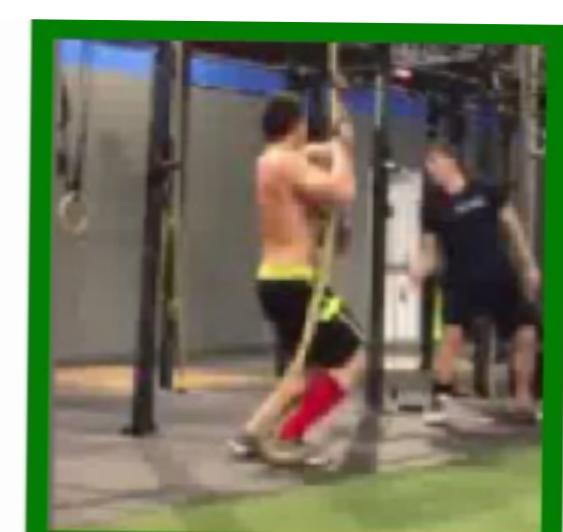
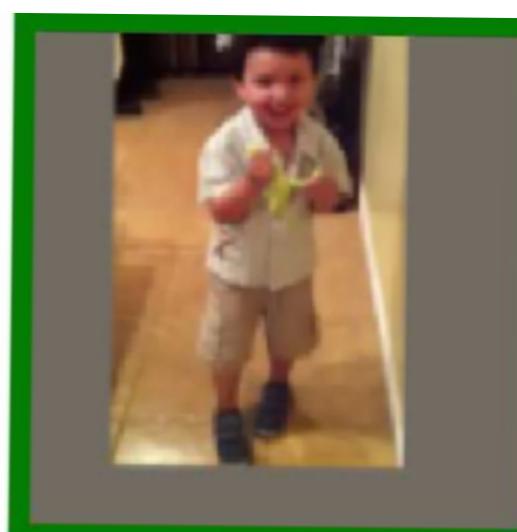


○ VGG-S ○ Chance ○ Mean Query ○ Artificial-Pb ○ Question-Aware Mask

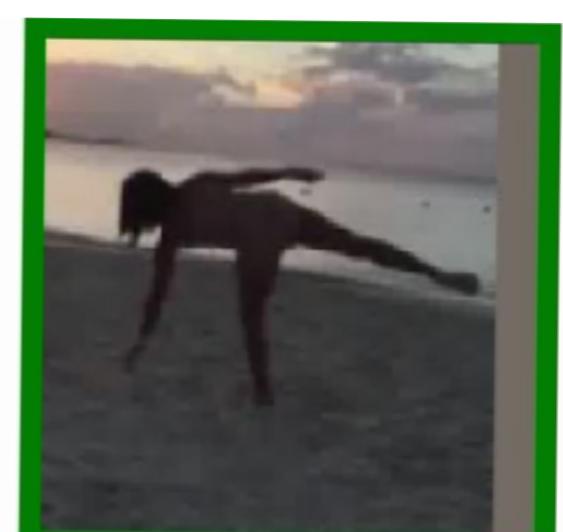
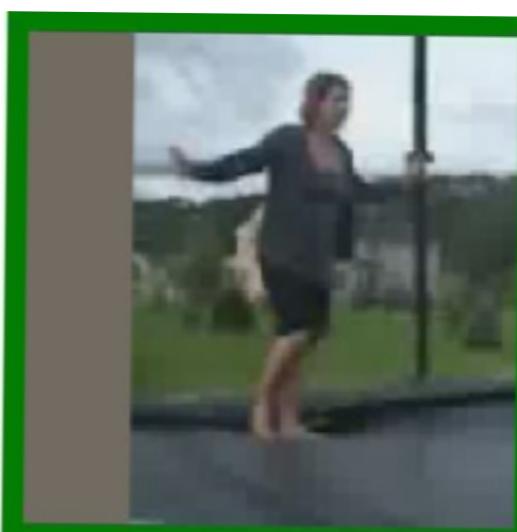
pelvis is bent



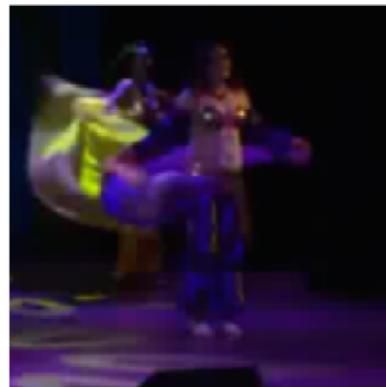
**right wrist
is near
neck**



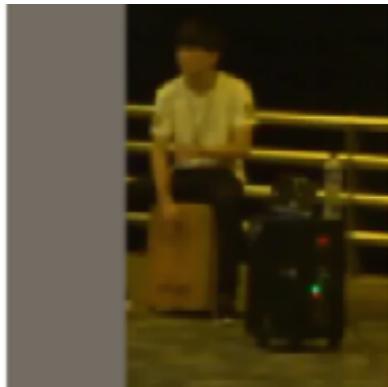
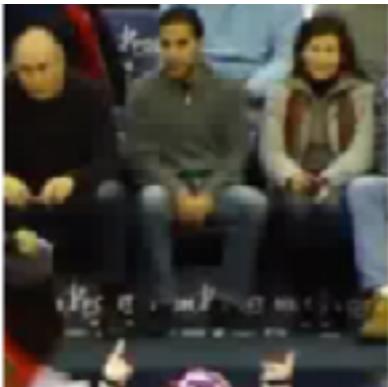
**left wrist
is beyond
left elbow**



pelvis is bent



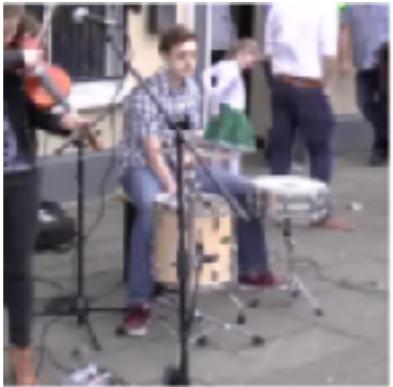
**right knee
is far from
left knee**



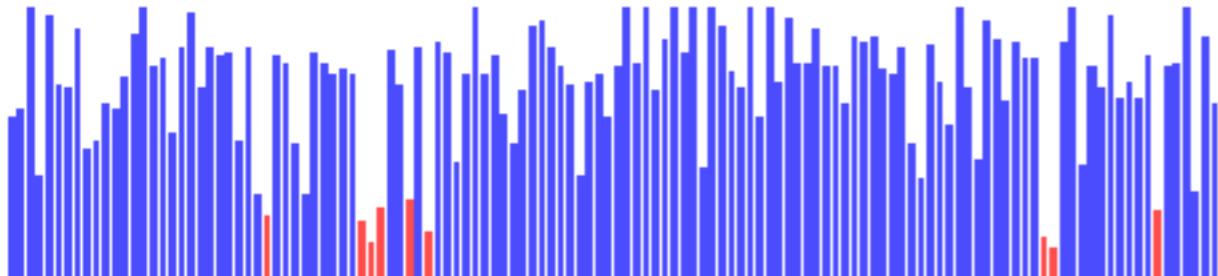
**left wrist
is beyond
neck**



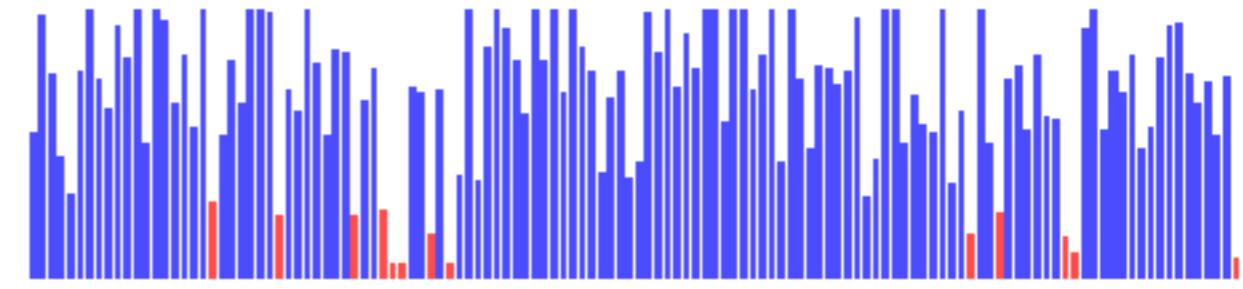
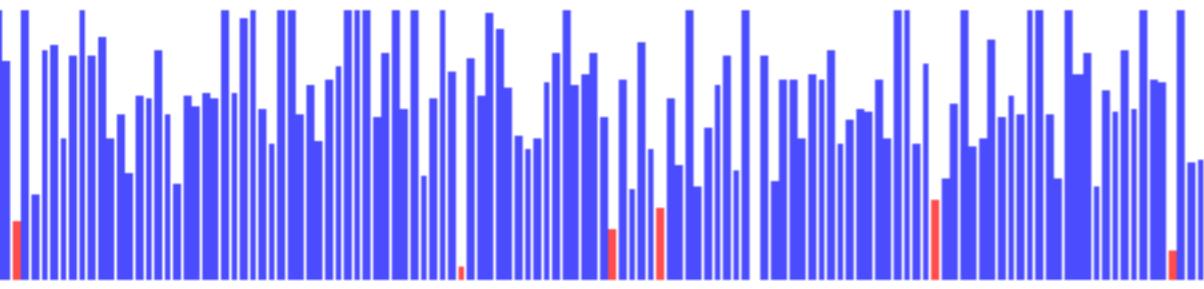
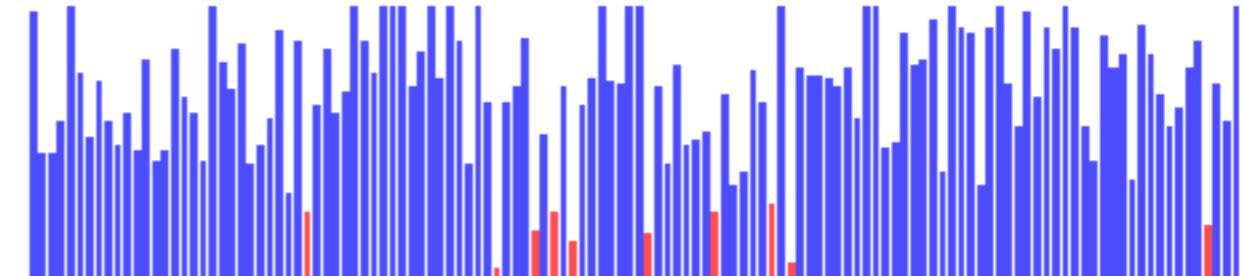
**right elbow
is near
left elbow**



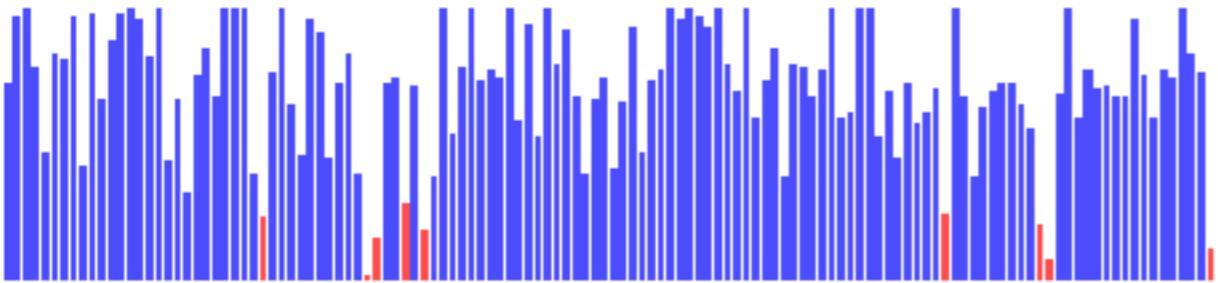
“right shoulder is (not) **bent**”



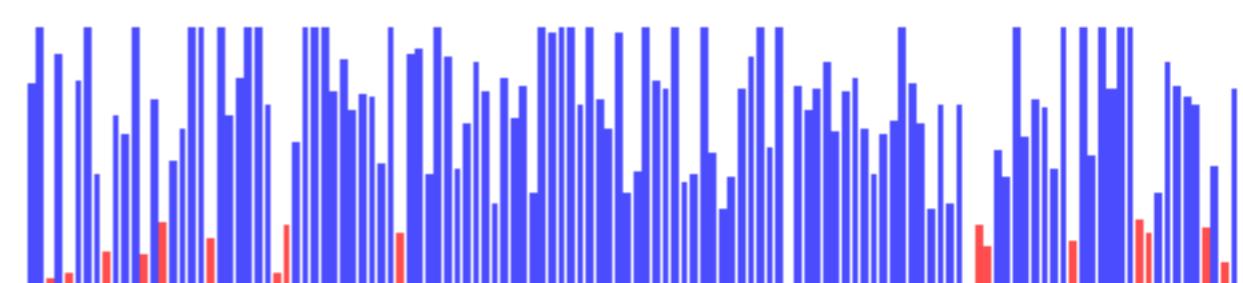
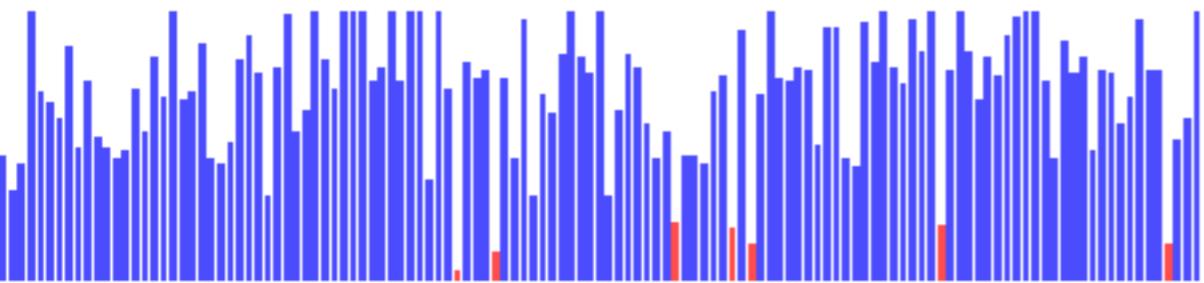
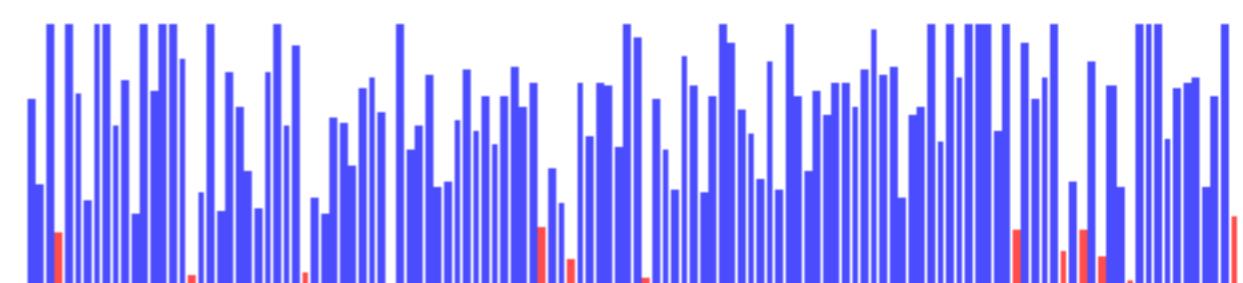
“left wrist is (not) **near** head”



“left wrist is (not) **beyond** pelvis”



“left knee is (not) **near** right knee”



Main Contributions

Pose-Aware embedding networks

Image-Language embedding network

Two solutions for language subsets

A dataset of 2D, 3D, and Language
Primitive descriptors