# COMBINING IMAGES AND WORDS IN DEEP NETWORKS THAT IDENTIFY PEOPLE FROM BODY SHAPE

## ALICE J. O'TOOLE
### *THE UNIVERSITY OF TEXAS AT DALLAS*

# ACKNOWLEDGEMENTS

- UT-Dallas
  - Blake A. Myers
  - Matthew Q. Hill
  - Veda Nandan Gandi
  - Thomas Metz
  - Lucas Jaggernauth
  - Madeline Rachow

- Johns Hopkins Univ.
  - Rama Chellappa
  - Carlos D. Castillo (Amazon)
  - Ram Prabhakar
- STR – team lead

- UT-Dallas
  - Matt Hill
  - Carina Hahn
- MPI for Intelligent Machinery
  - Stephan Streuber
  - Michael Black

# OVERVIEW

- Problem – person identification based on body shape
  - Biometric Recognition and Identification at Altitude and Range (BRIAR)
    - IARPA https://www.iarpa.gov/research-programs/briar

- Linguistic descriptors to "quantify" body shape
  - psychology, computer graphics

- Body identification networks:
  - linguistic descriptors
  - object-based shape descriptors

- Person recognition = face + body + gait
  - fusion
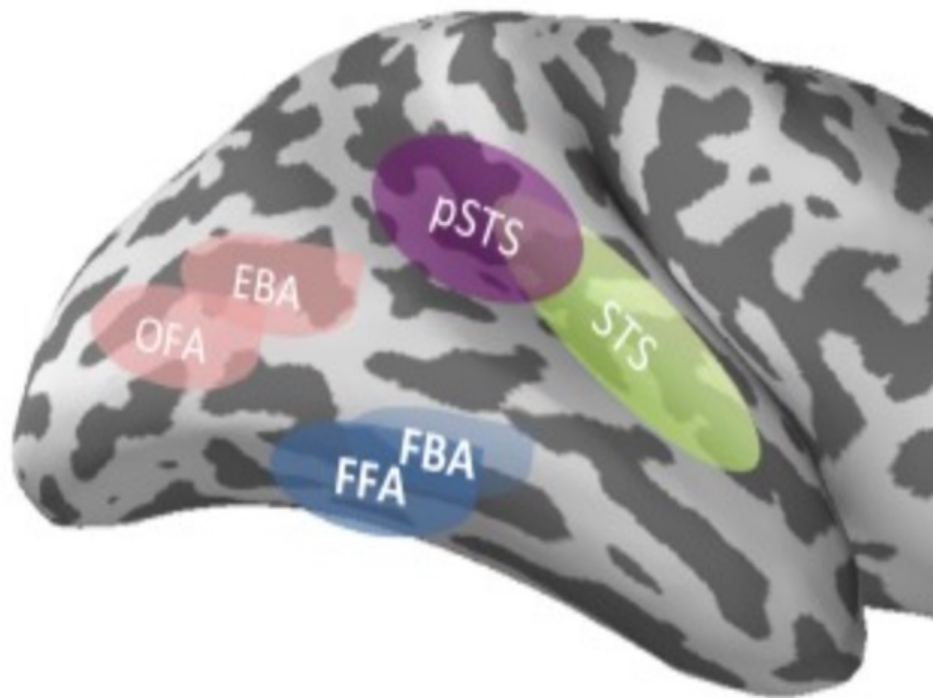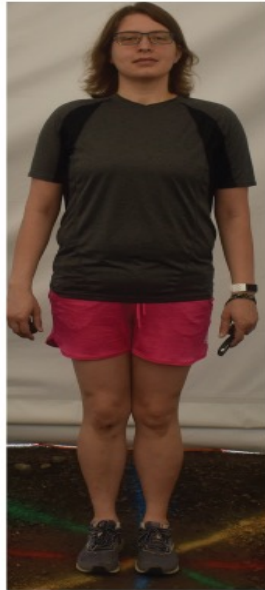
3

# PROBLEM



*face*

*body*

*gait*

4

Yovel & O'Toole (2016)

controlled    close range    UAV

100m    200m    400m    500m

subject consented to publication

*same person or different people?*



face

*body*

gait

7

# BODY AS A BIOMETRIC

- Why use body?
  - visible at large distances
  - subset of cues constant over change in view
    - height, weight, proportions, rough shape
  - "fusability"

  - *You have no other option!*

# BODY IS LEAST COMMON DENOMINATOR

subjects consented to publication

# BODY AS A BIOMETRIC

- Why not use body?
  - not unique
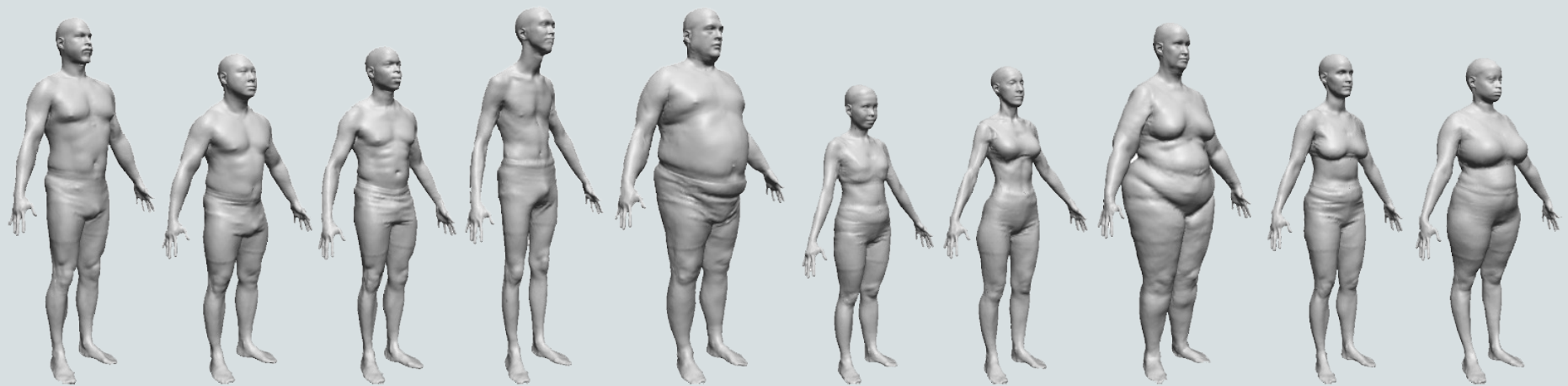  - lack of a body algorithms
  - •face
  - gait
    - body



**Bo**...**n**

# LINGUISTIC DESCRIPTIONS
# &
# 3D BODY SHAPES

Hill, Matthew Q., et al. "Creating body shapes from verbal descriptions by linking similarity spaces." *Psychological science* 27.11 (2016): 1486-1497.

Streuber, Stephan, et al. "Body talk: Crowdshaping realistic 3D avatars with words." *ACM Transactions on Graphics (TOG)* 35.4 (2016): 1-14.

# HUMAN BODY SHAPE



- body = complex 3D shape
  - Laser scan = 12500 vertices and 25000 facets

# LINGUISTIC DESCRIPTIONS OF BODIES

- muscular, athletic

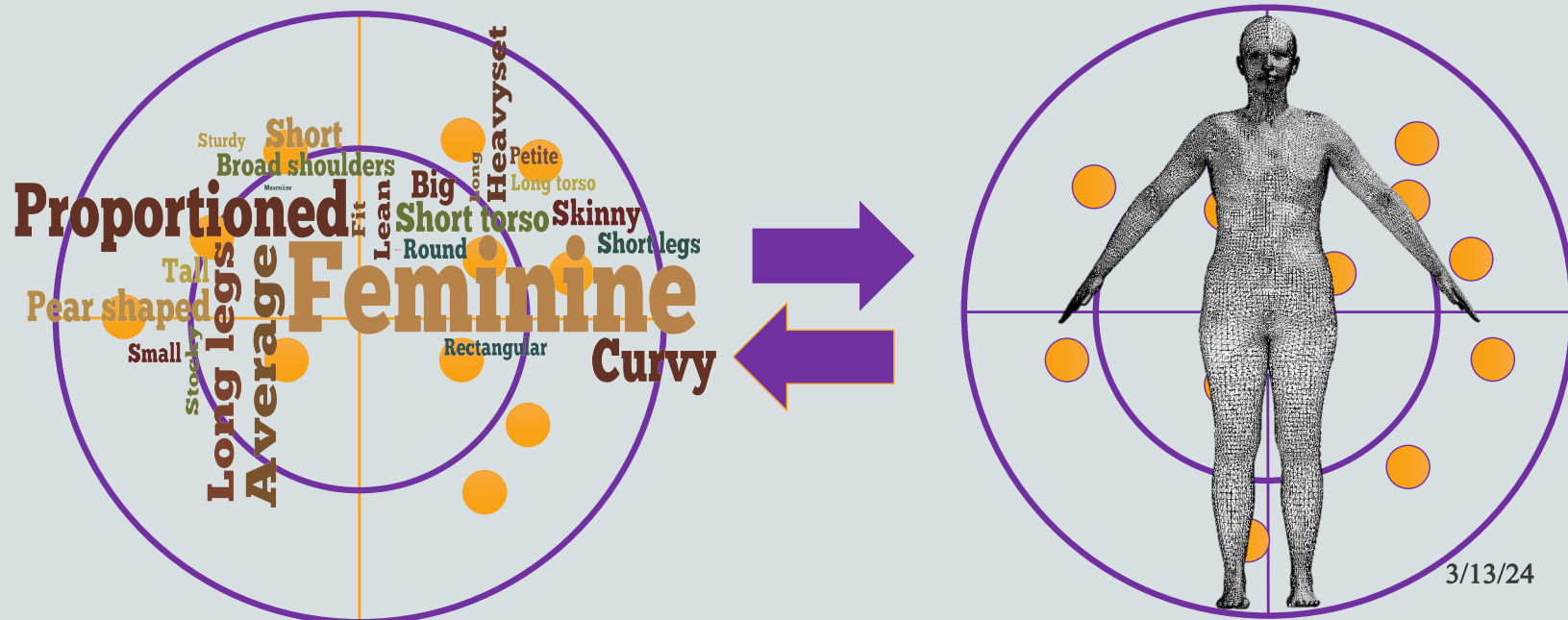- stout, portly

- shapely, hourglass

# RATIONALE

- human language and vision
  - evolved a long time ago
    - 50K and 2 million years ago
    - words don't leave fossils or tool fragments
  - language communicates information efficiently
    - not every vertex in the laser scan carries a lot of information

long legs     lean     curvy

# APPROACH

- human descriptions to create a similarity space

  - point proximity -> similarity between *body descriptions*

- geometric shape space to *ground-truth* description space

  - point proximity -> similarity between *body shapes*

# BODY DESCRIPTIONS



|  | Does Not Apply | Applies Somewhat | Applies Perfectly |
|---|---|---|---|
| Proportioned | ● | ○ | ○ |
| Rectangular | ● | ○ | ○ |
| Stocky | ● | ○ | ○ |
| Short legs | ● | ○ | ○ |
| Muscular | ● | ○ | ○ |
| Average | ● | ○ | ○ |
| Tall | ● | ○ | ○ |
| Sturdy | ● | ○ | ○ |
| Big | ● | ○ | ○ |
| Long legs | ● | ○ | ○ |
| Lean | ● | ○ | ○ |
| Short torso | ● | ○ | ○ |
| Pear shaped | ● | ○ | ○ |
| Petite | ● | ○ | ○ |
| Broad shoulders | ● | ○ | ○ |
| Heavyset | ● | ○ | ○ |
| Long | ● | ○ | ○ |
| Long torso | ● | ○ | ○ |
| Round (Apple) | ● | ○ | ○ |
| Built | ● | ○ | ○ |
| Fit | ● | ○ | ○ |
| Skinny | ● | ○ | ○ |
| Masculine | ● | ○ | ○ |
| Small | ● | ○ | ○ |
| Short | ● | ○ | ○ |
| Feminine | ● | ○ | ○ |
| Curvy | ● | ○ | ○ |

NEXT PAGE

3/13/24

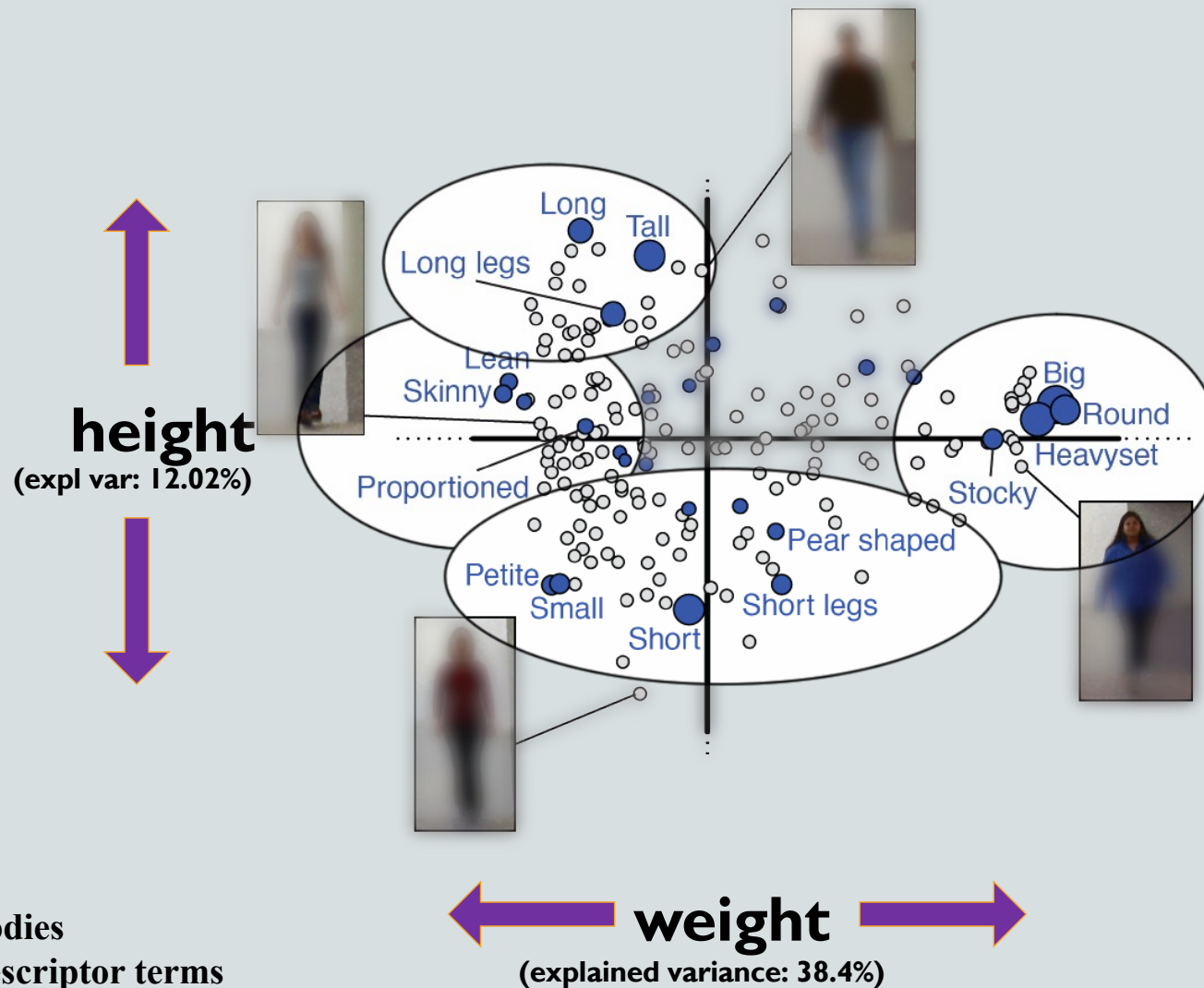| Descriptor Terms | | | |
|---|---|---|---|
| proportioned | sturdy | broad shoulders | skinny |
| rectangular | big | heavyset | masculine |
| stocky | long legs | long | small |
| short legs | lean | long torso | short |
| muscular | short torso | round (apple) | feminine |
| average | pear shaped | built | curvy |
| tall | petite | fit | |

3/13/24

# LANGUAGE SPACE DATA

- Body representations:
  - descriptions made from <u>images</u> of people

3/13/24

# LANGUAGE SIMILARITY SPACE

- applied **correspondence analysis** to:
  - descriptor vectors for the 164 female bodies
    - 27 elements - terms that "applied perfectly" to the body

- Correspondence Analysis (CA) (Benzicri, 1973)
  - multivariate analysis analogous to PCA, but for categorical data
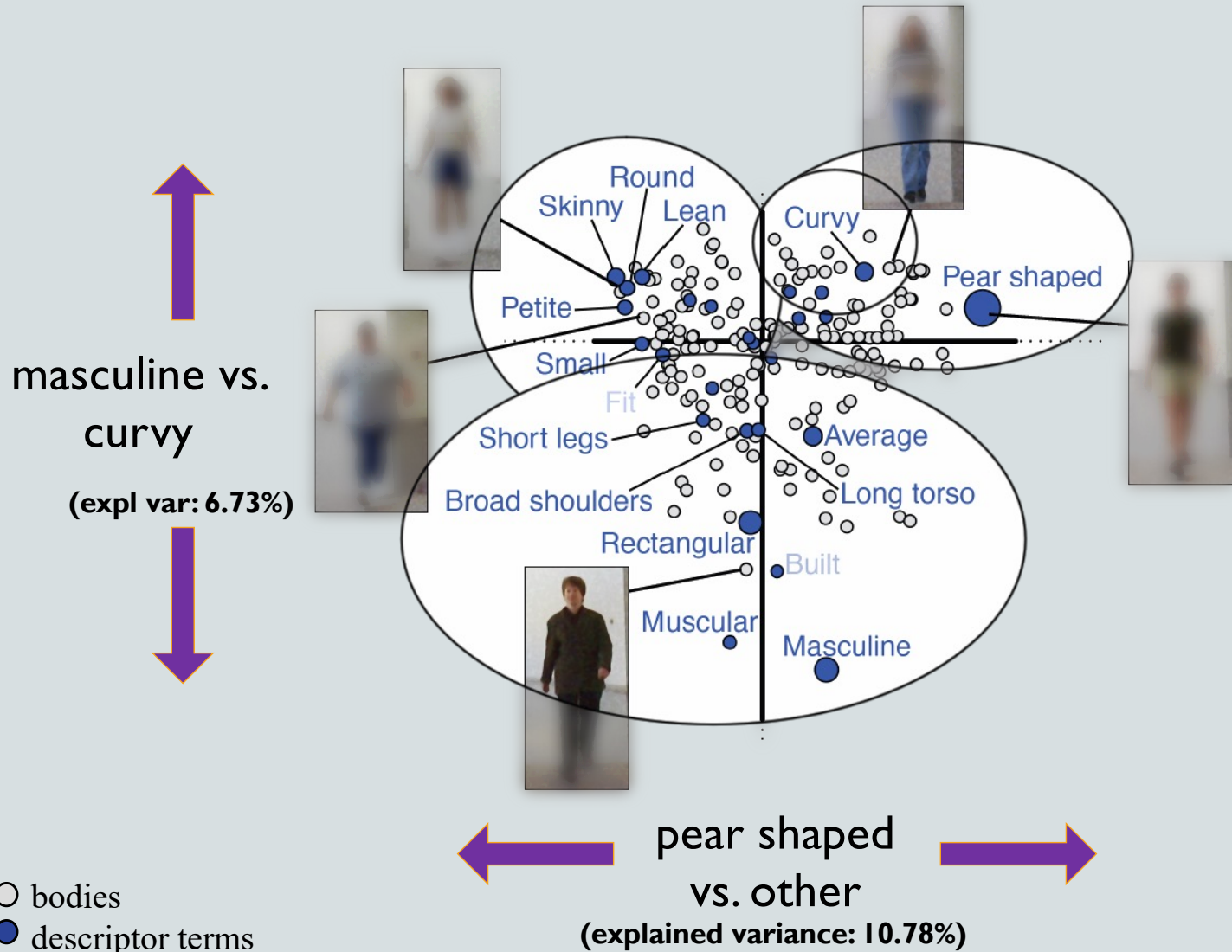  - allows observations (bodies) to be plotted in the same space as the variables (descriptor terms)

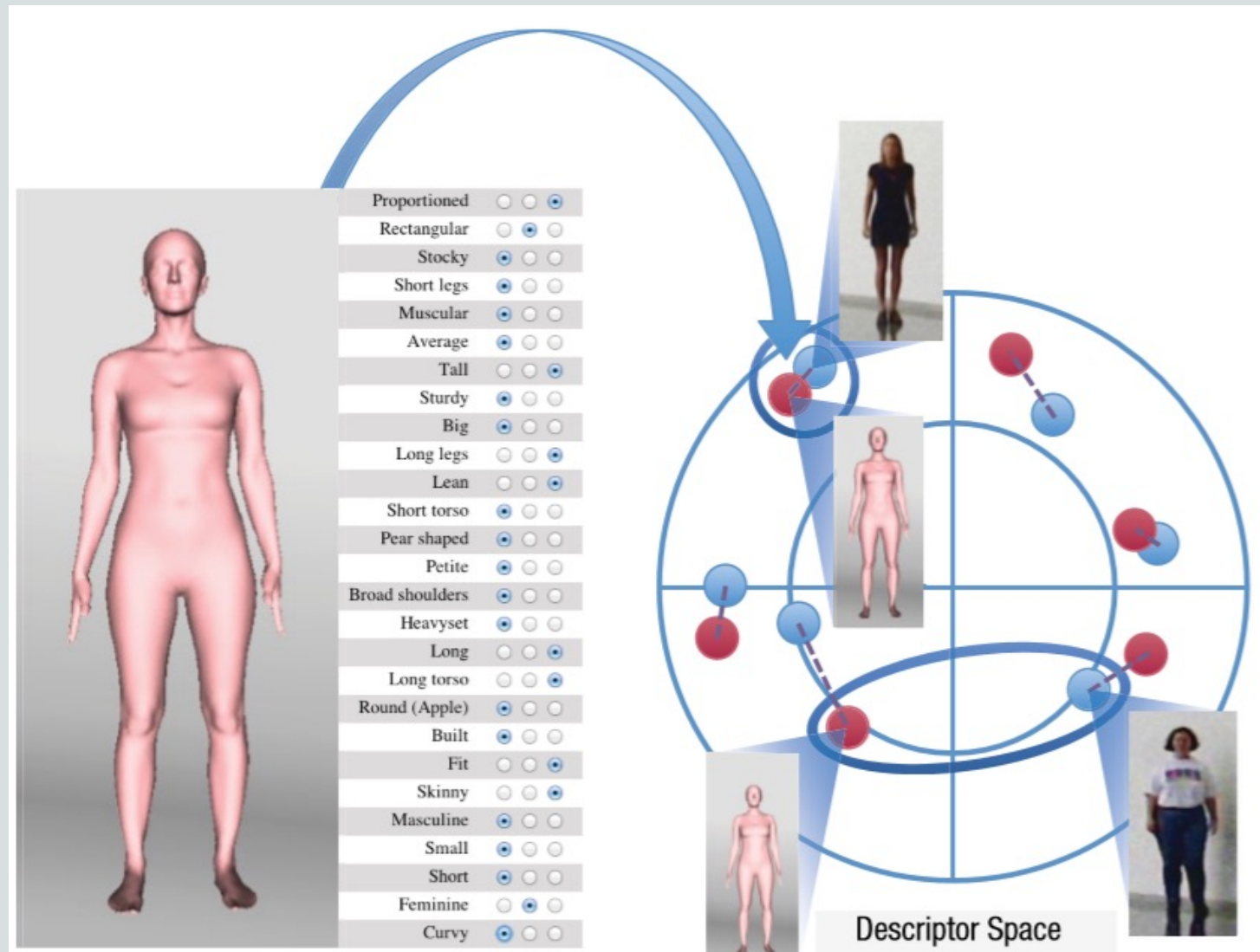3/13/24

# LANGUAGE SPACE

## (1ST & 2ND AXES)



**height**
(expl var: 12.02%)

Long
Tall
Long legs
Lean
Skinny
Proportioned
Petite
Small
Short
Short legs
Pear shaped
Big
Round
Heavyset
Stocky

**weight**
(explained variance: 38.4%)

○ bodies
● descriptor terms

3/13/24

# LANGUAGE SPACE

## (3ST & 4TH AXES)



masculine vs. curvy

**(expl var: 6.73%)**

Round
Skinny  Lean  Curvy
Pear shaped
Petite
Small
Fit
Short legs  Average
Broad shoulders  Long torso
Rectangular  Built
Muscular
Masculine

pear shaped vs. other
**(explained variance: 10.78%)**

○ bodies
● descriptor terms

3/13/24

Body model PCA of laser scans of bodies (Loper et al., 2016)
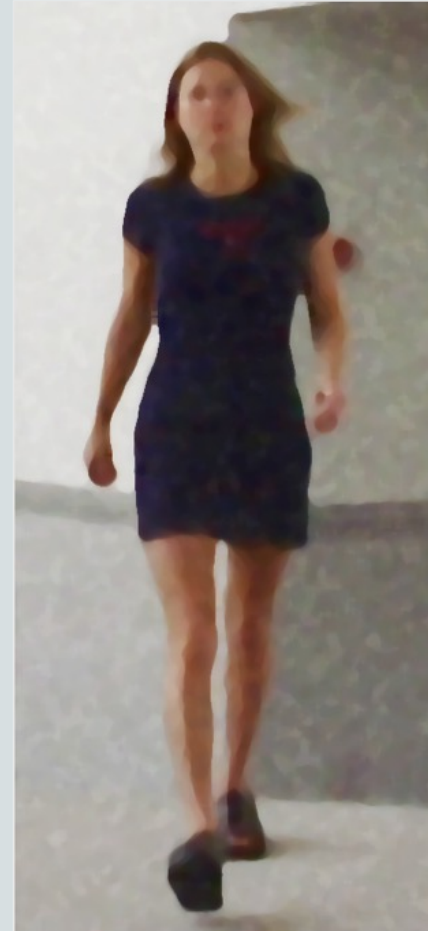
23

# 3D BODY SYNTHESIS FROM DESCRIPTIONS



descriptions

PCA coefficients

PCA of 3000+ laser scans
SMPL model (Loper et al., 2016)

24

subject consented to publication
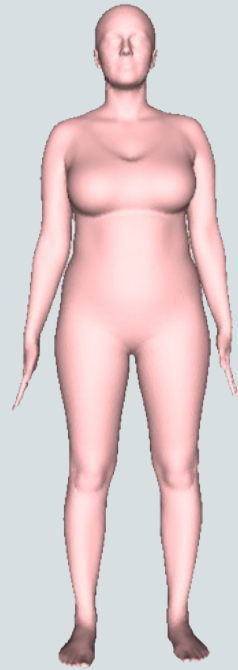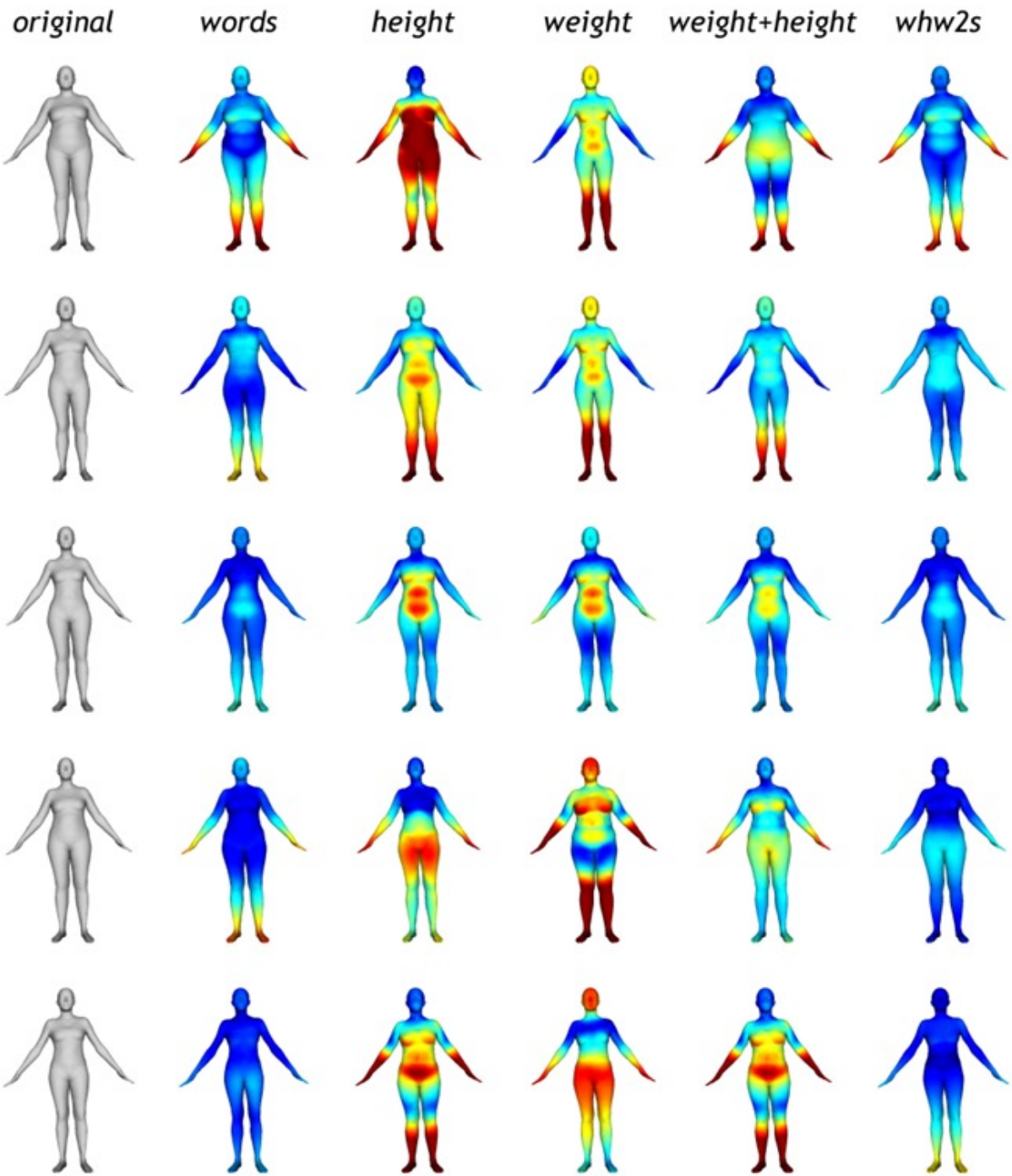
subject consented to publication

3/13/24

subject consented to publication

3/13/24

original   words   height   weight   weight+height   whw2s

Streuber et al. (2016)

29

# CONCLUSIONS

- linguistic descriptions
  - can be used to synthesize 3D bodies

- efficient way to perform a laser scan without a laser scanner ☺
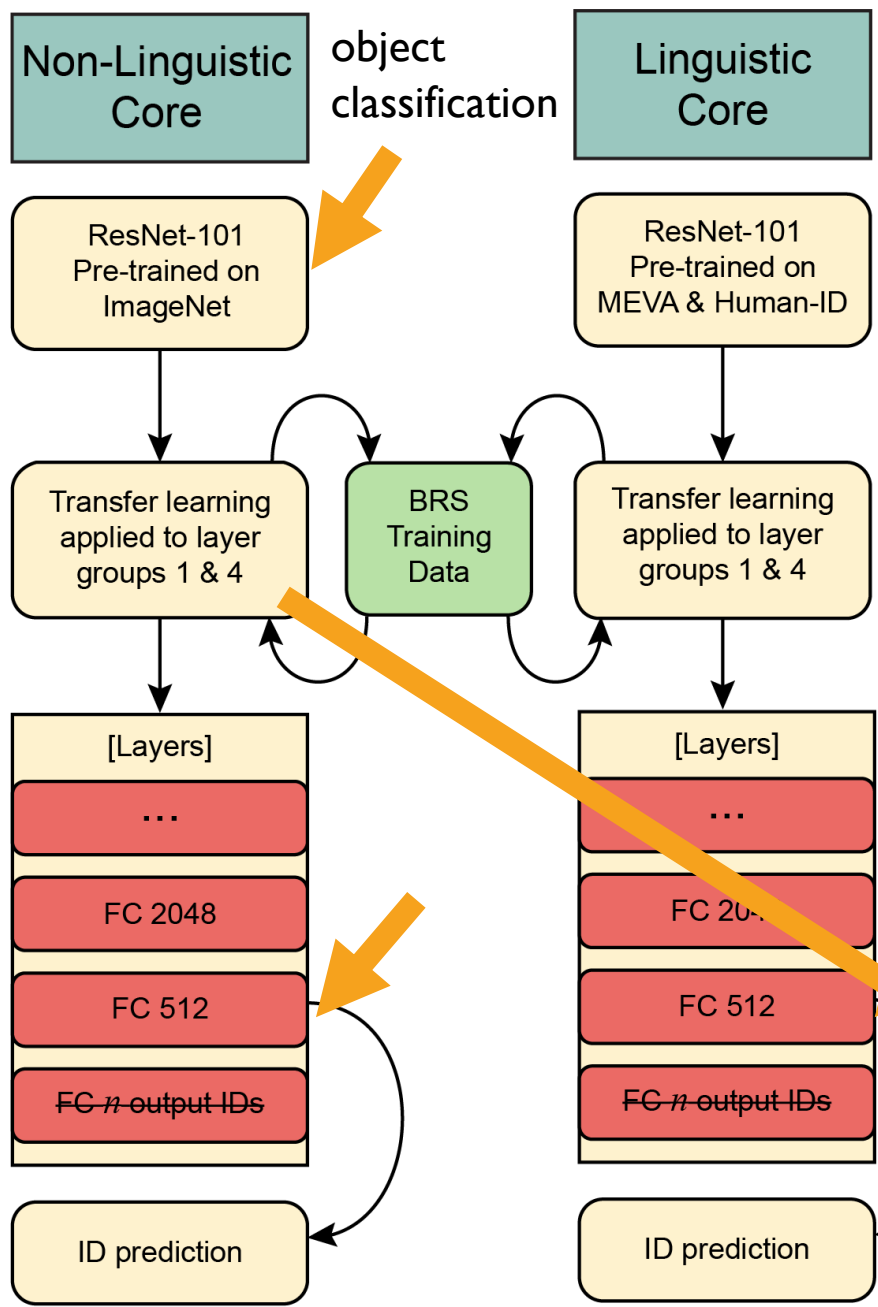
# IDENTIFICATION FROM BODY SHAPE

Myers BA, Jaggernauth L, Metz TM, Hill MQ, Gandi VN, Castillo CD, O'Toole AJ. Recognizing People by Body Shape Using Deep Networks of Images and Words. arXiv:2305.19160. 2023 May 30. *Proc. IEEE International Joint Conference on Biometric, Sept. 2023*

# WORDS FOR BODY IDENTIFICATION

- Rationale
  - **descriptors sufficient to synthesize 3D body**
  - descriptor-based representation for identification?

- Advantages
  - robust across large distances
  - generalize across yaw and pitch (curvy, tall, stout, long legs,)
    - accessible across a range of view
  - (relatively) clothing independent

  - Explainable AI??

## CURRENT PROBLEM

- learn mapping from images to descriptors
  - pretraining – to categorize body shape


- image to identity
  - transfer learning – image to identity
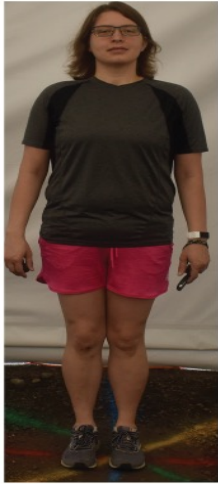  - fine tuning within a category

Non-Linguistic Core

object classification

Linguistic Core

ResNet-101 Pre-trained on ImageNet

ResNet-101 Pre-trained on MEVA & Human-ID

Transfer learning applied to layer groups 1 & 4

BRS Training Data

Transfer learning applied to layer groups 1 & 4

[Layers]

...

FC 2048

FC 512

FC *n* output IDs

ID prediction

[Layers]

...

FC 2048

FC 512

FC *n* output IDs

ID prediction

curvy,
tall,
stocky,
short legs
muscular,.....etc.

**identity trained**
close range,
UAV
100m, 200m, etc....1000m

**identity trained**
close range,
UAV
100m, 200m, etc....1000m

34

# MODELS

- linguistic body model (LCRIM)
  - linguistic core model
    - body image to linguistic description
  - identity-tuning
    - body image to identity

- non-linguistic body model (NLCRIM)
  - pre-trained object classification core model
    - ImageNet trained
  - identity-tuning
    - body image to identity

- Fusion = LCRIM + NLCRIM

controlled    close range    UAV

100m    200m    400m    500m

- **training**
  - 577 IDs
    - 242,386 images

- **test**
  - 485 gallery IDs
    - 43,722 images
  - 260 probe IDs
    - 2,192,305 image frames from 9,795 videos
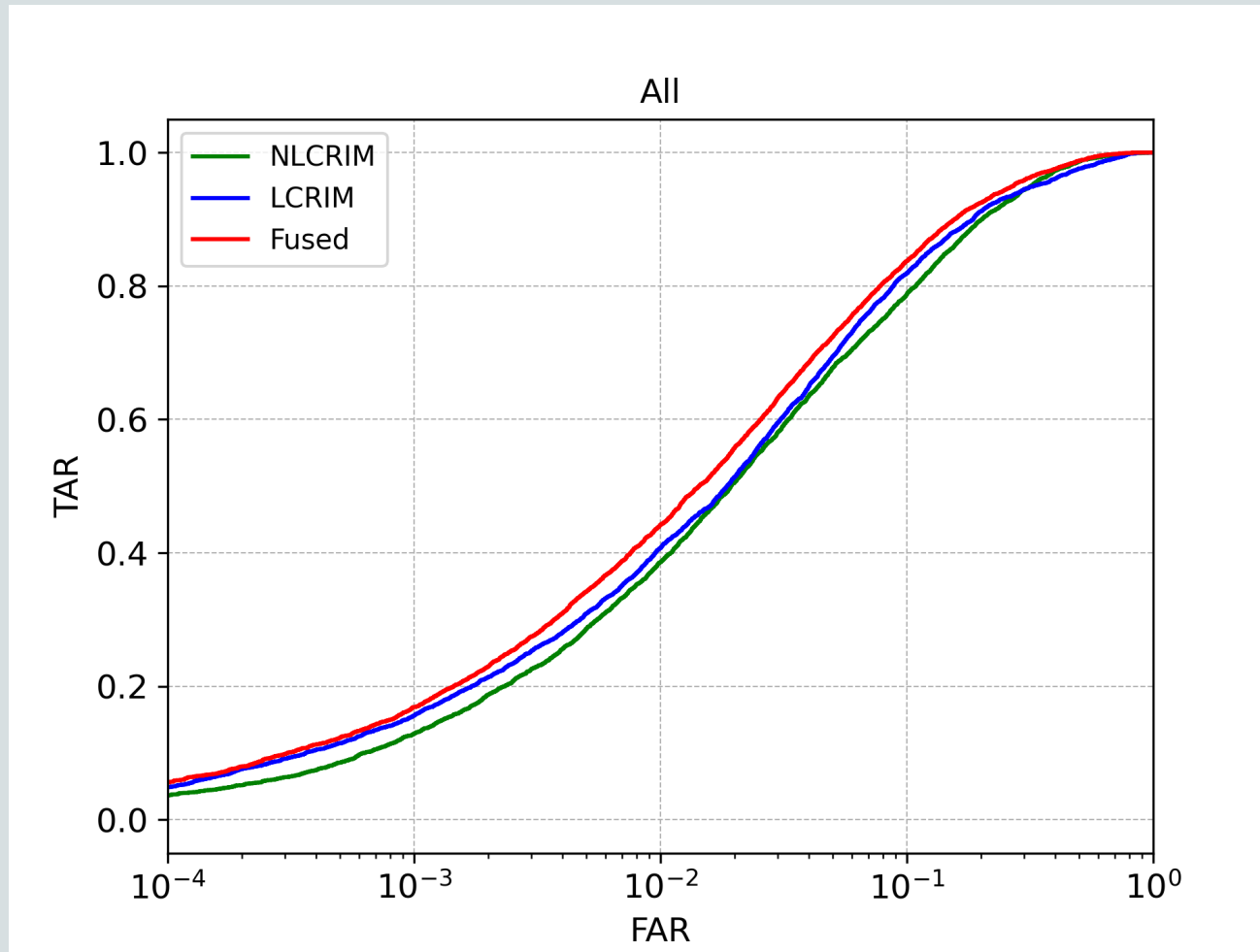
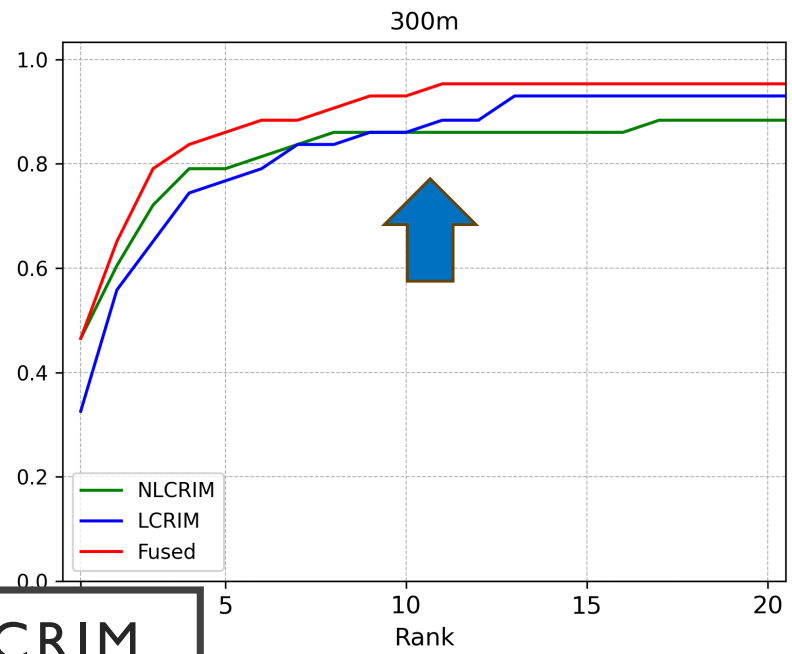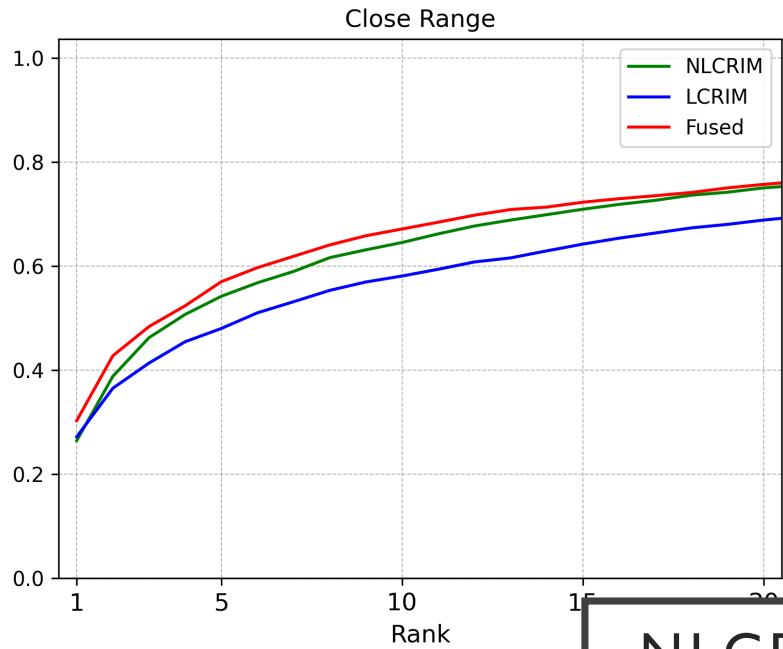(BRS-BTS dataset, Cornett, et. al., 2022)

# Cumulative Match Characteristic
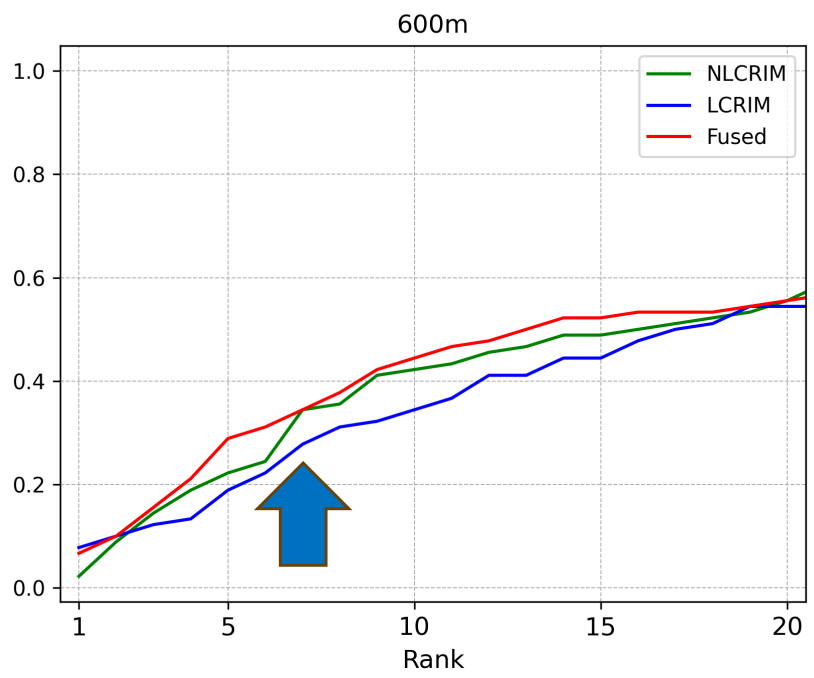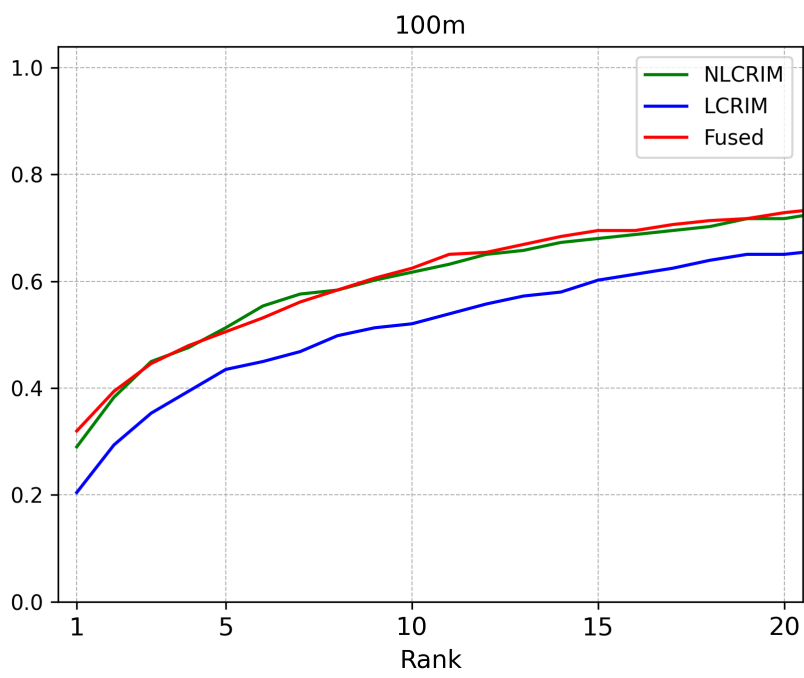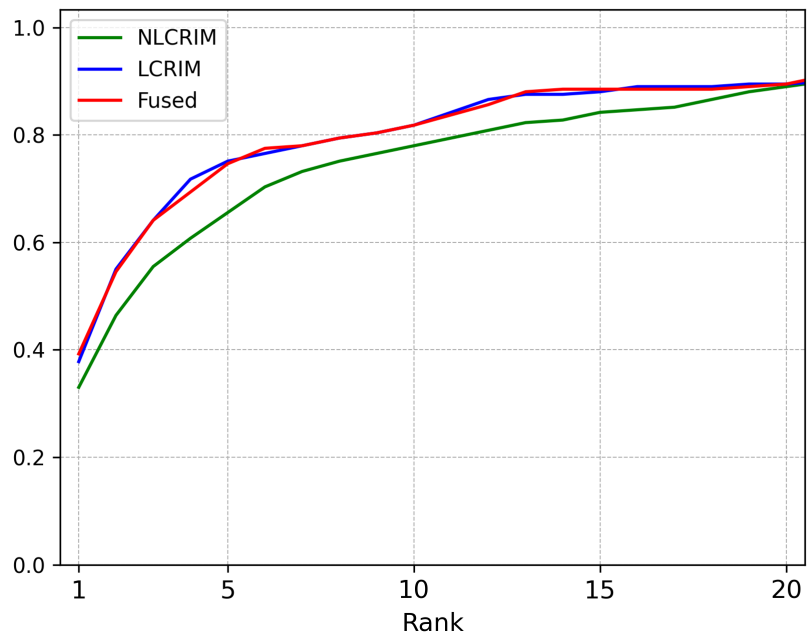
# Receiver Operating Characteristic Curve



All

# DISTANCE CONDITIONS

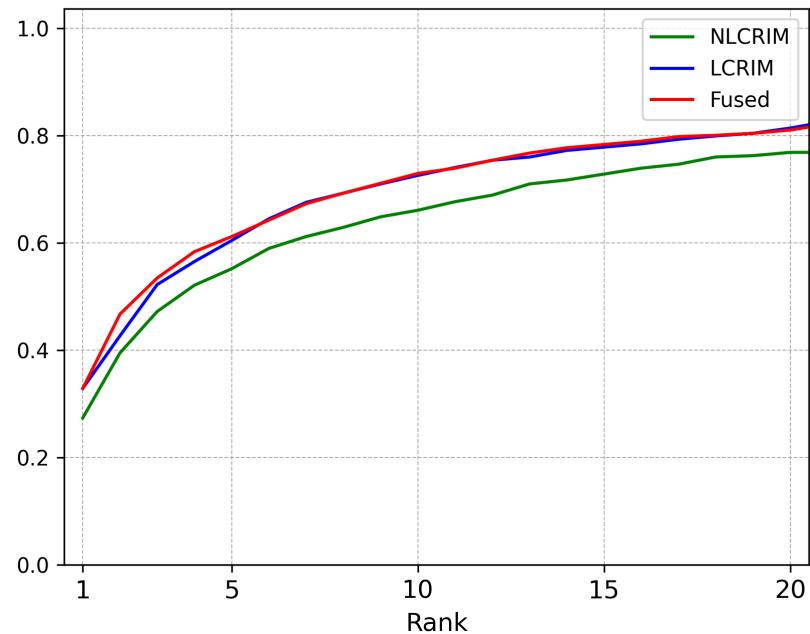- Linguistic > as views and pitch get more extreme?
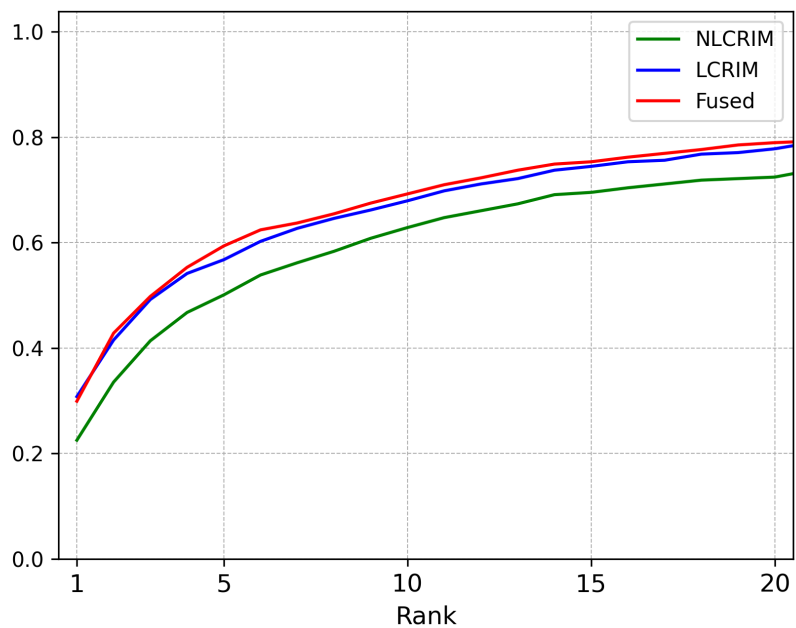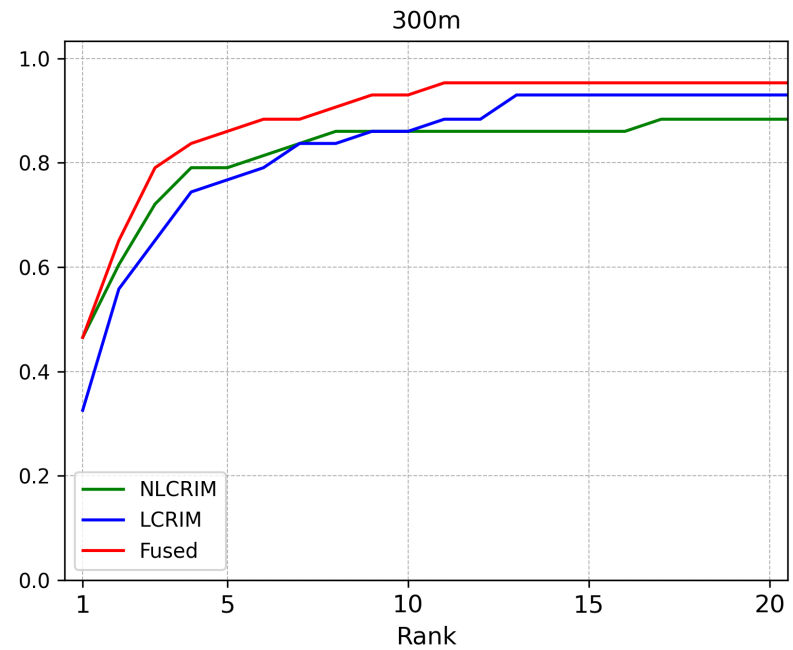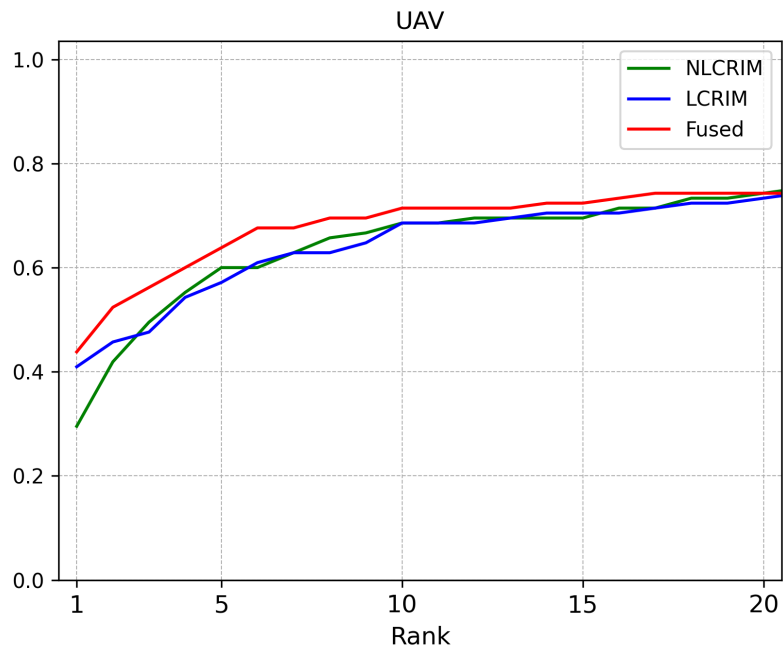
Close Range

300m

NLCRIM > LCRIM

100m

600m

LCRIM > NLCRIM

42

# FUSION >> (NLCRIM OR LCRIM)

## LINGUISTIC? NON-LINGUISTIC? FUSED?

- condition-dependent
  - fusion *almost* always best


  - linguistic/non-linguistic
    - less predictable

# CONCLUSIONS

- Linguistic descriptors
  - complement body shape representations
  - better at further distances (tentatively)
  - tap similar types of information

# PERSON = FACE + BODY + GAIT

FUSION, VARIANCE, QUALITY

ID₁     ID₂

LIMITS OF THE BODY

same person or different people?

same person or different people?

same person or different people?

47

subject consented to publication

# FACE, BODY, & GAIT: MODEL (DIS)AGREEMENT

|  | | Body 1 | Body 2 | Body 3 | Face 1 | Face 2 | Gait |
|---|---|---|---|---|---|---|---|
| | | 0 | 1 | 2 | 3 | 4 | 5 |
| Body 1 | 0 | 1 | 0.44643 | 0.446626 | 0.186319 | 0.185375 | 0.15463 |
| Body 2 | 1 | 0.44643 | 1 | 0.785518 | 0.25484 | 0.255659 | 0.346526 |
| Body 3 | 2 | 0.446626 | 0.785518 | 1 | 0.250673 | 0.251236 | 0.352391 |
| Face 1 | 3 | 0.186319 | 0.25484 | 0.250673 | 1 | 0.660499 | 0.127177 |
| Face 2 | 4 | 0.185375 | 0.255659 | 0.251236 | 0.660499 | 1 | 0.132766 |
| Gait | 5 | 0.15463 | 0.346526 | 0.352391 | 0.127177 | 0.132766 | 1 |

# Fusion

# Face Fusion

# Body Fusion



CMC Protocol 4.1sum_fusion

CMC Protocol 4.1face_fusion

CMC Protocol 4.1body_fusion

Receiver Operating Characteristic (ROC) Curve: sum_fusion — ROC curve (area = 0.97)

Receiver Operating Characteristic (ROC) Curve: face_fusion — ROC curve (area = 0.91)

Receiver Operating Characteristic (ROC) Curve: body_fusion — ROC curve (area = 0.95)

Yovel & O'Toole (2016)

# APPROACH

- fusion on a case-by-case basis
  - requires **quality** of face vs. body vs. gait with limited meta-data

- What happens when they do not agree?
  - face with body?
  - face with face? body with body?
  - gait with face or body?

- *Can disagreement be informative of quality???*

# VARIANCE OF ESTIMATES

- Proposal
  - Can variance of model estimates guide fusion?

- Predict
  - high variance indicates "low quality" and low accuracy

- Prerequisite (sanity test)
  - *Does variance of model estimates relate to accuracy?*

# DOES MODEL VARIANCE PREDICT ACCURACY?

- ## Variance *on each item*
  - all-model variance
  - face-model variance
  - body-model variance

- ## Performance:
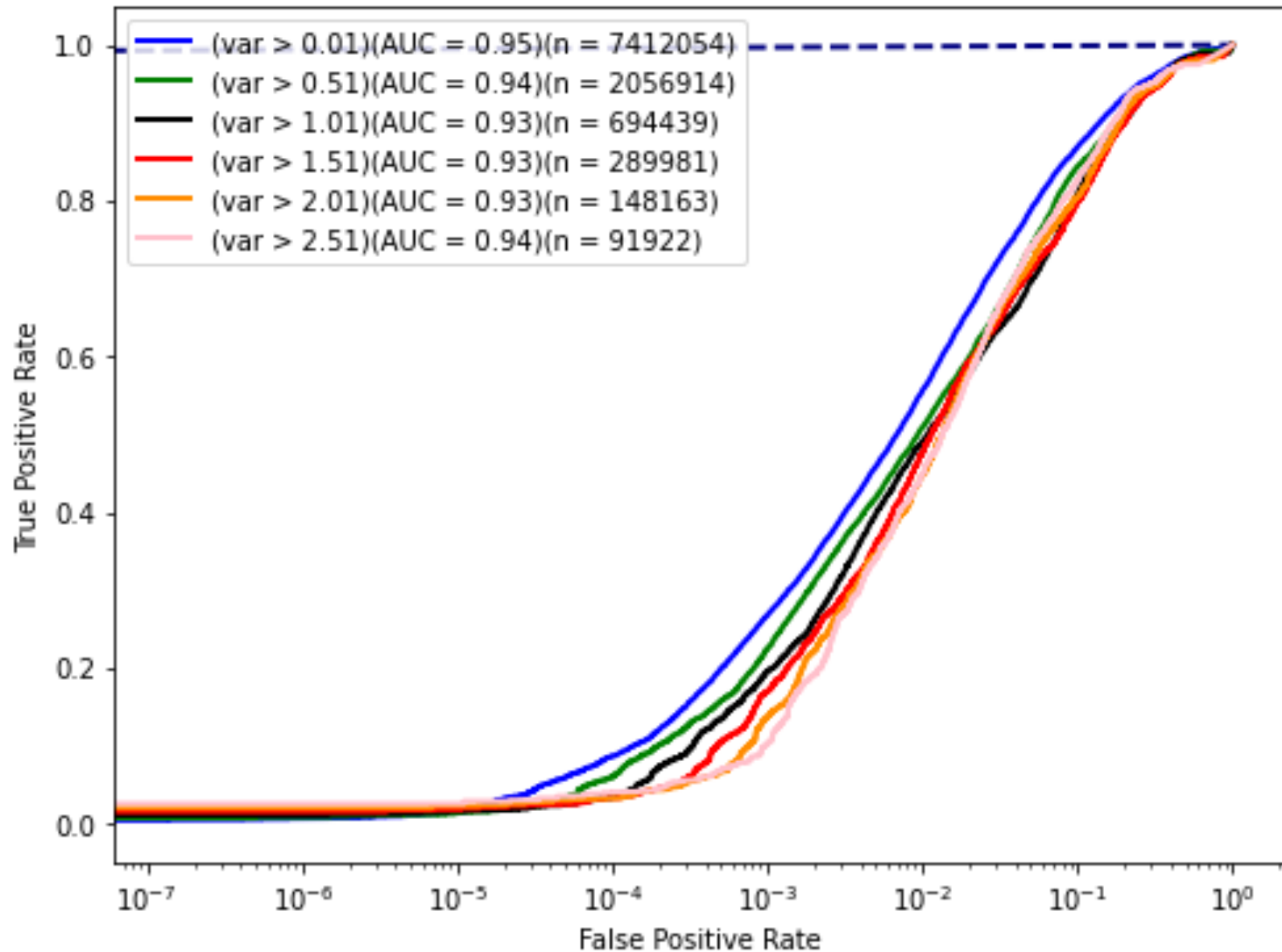  - face fusion similarity scores (2 face algorithms)
  - body fusion similarity scores (3 body algorithms)

Receiver Operating Characteristic (ROC) Curve: - All Variance Body Fusion

Legend:
- (var > 0.01)(AUC = 0.95)(n = 7615206)
- (var > 0.76)(AUC = 0.93)(n = 2133375)
- (var > 1.51)(AUC = 0.90)(n = 510945)
- (var > 2.26)(AUC = 0.90)(n = 230017)
- (var > 3.01)(AUC = 0.91)(n = 141567)
- (var > 3.76)(AUC = 0.92)(n = 97844)

Low variability model scores - better performance with body information

Low variability body models better performance with body

Receiver Operating Characteristic (ROC) Curve: - Face Variance Body Fusion

Legend:
- (var > 0.01)(AUC = 0.95)(n = 6267408)
- (var > 0.26)(AUC = 0.95)(n = 2104272)
- (var > 0.51)(AUC = 0.94)(n = 1043662)
- (var > 0.76)(AUC = 0.93)(n = 578747)
- (var > 1.01)(AUC = 0.93)(n = 341581)
- (var > 1.26)(AUC = 0.92)(n = 211138)

Low variability face models better performance with body

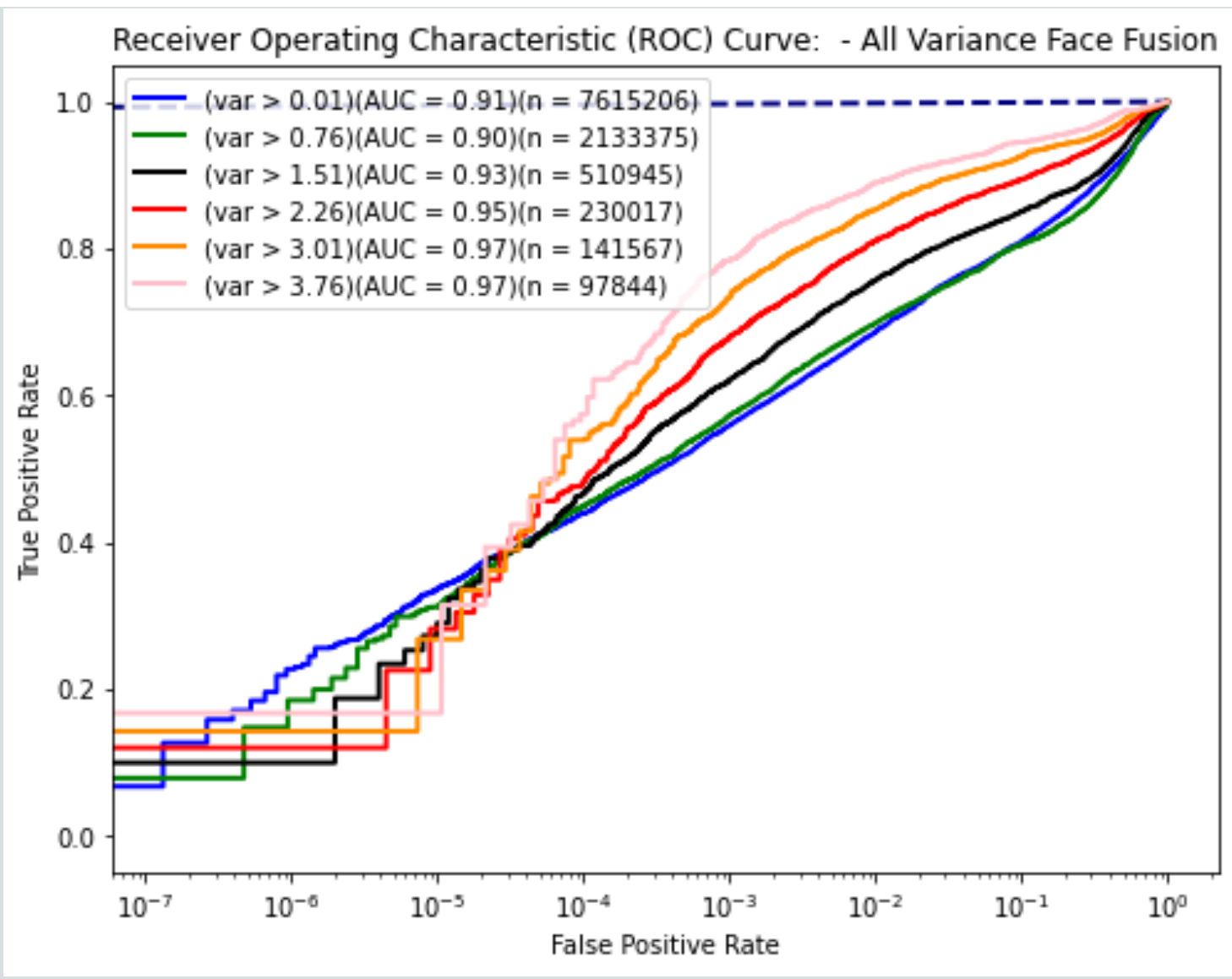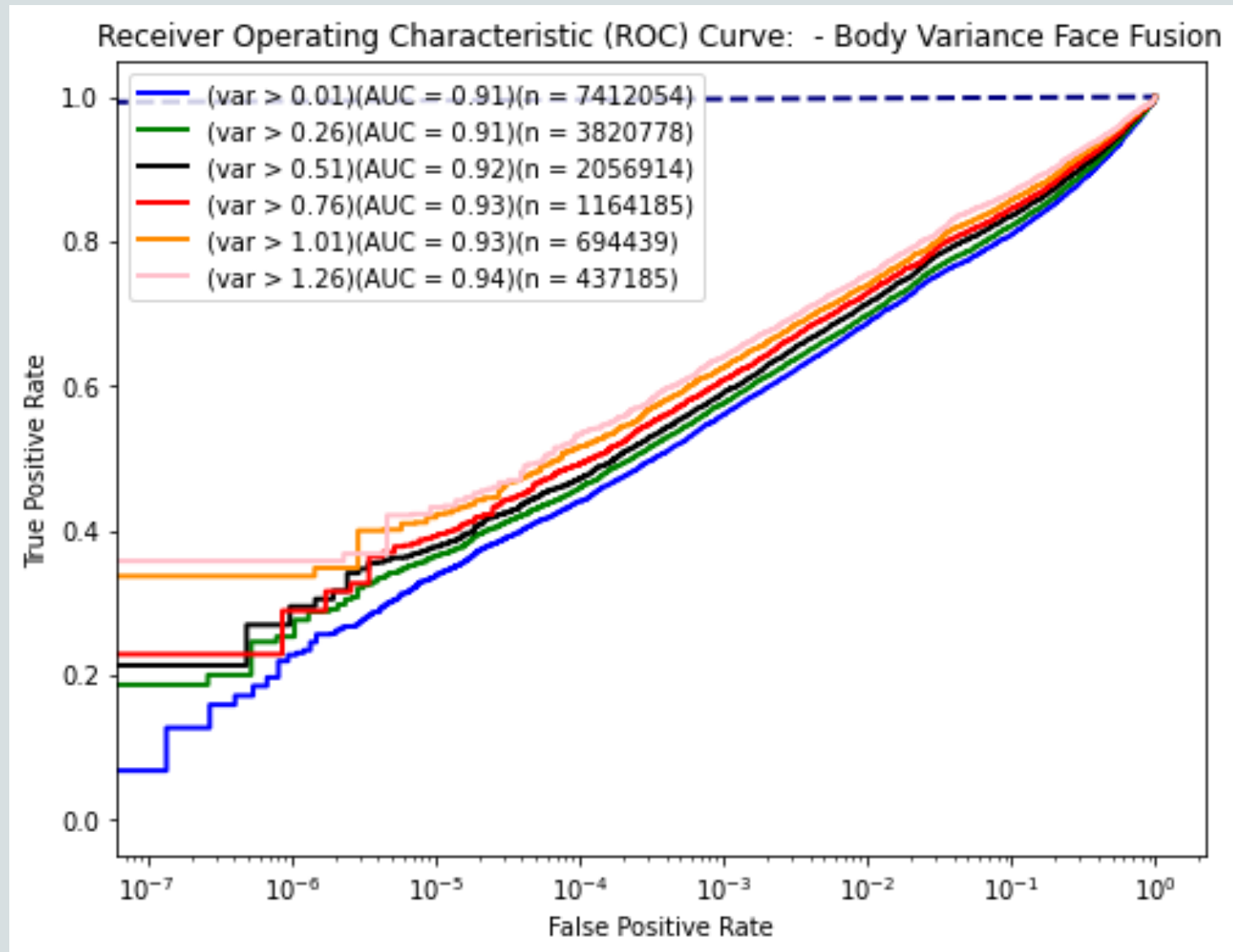Receiver Operating Characteristic (ROC) Curve: - Face Variance Face Fusion

- (var > 0.01)(AUC = 0.91)(n = 6267408)
- (var > 0.26)(AUC = 0.91)(n = 2104272)
- (var > 0.51)(AUC = 0.92)(n = 1043662)
- (var > 0.76)(AUC = 0.92)(n = 578747)
- (var > 1.01)(AUC = 0.92)(n = 341581)
- (var > 1.26)(AUC = 0.93)(n = 211138)

Low variability face models better performance with the face

# 2 INVERSIONS

Receiver Operating Characteristic (ROC) Curve:  - All Variance Face Fusion

- (var > 0.01)(AUC = 0.91)(n = 7615206)
- (var > 0.76)(AUC = 0.90)(n = 2133375)
- (var > 1.51)(AUC = 0.93)(n = 510945)
- (var > 2.26)(AUC = 0.95)(n = 230017)
- (var > 3.01)(AUC = 0.97)(n = 141567)
- (var > 3.76)(AUC = 0.97)(n = 97844)

High variability model scores *better* performance with the face

(except at very low FP)

High variability body estimates better performance with face!

# TAKE HOME MESSAGE

- Biometrics has ignored the body on the (correct) premise that it is not "unique"
  - not unique ≠ not helpful

- Linguistic descriptions of bodies
  - graphics, shape classification, identification

- Body algorithms boost identification over
  - face
  - gait

- Quality estimates from model discord within/across modalities

# ACKNOWLEDGEMENTS

# THANK YOU!