

Judgment Studies

Lecture 03 (Modelling II)

Alessandro Vinciarelli



University
of Glasgow



Social AI



Engineering and
Physical Sciences
Research Council

Outline

- Judgment Studies
- Example: Personality Perception in AAR
- Conclusions

Outline

- Judgment Studies
- Example: Personality Perception in AAR
- Conclusions

The Judgment Studies

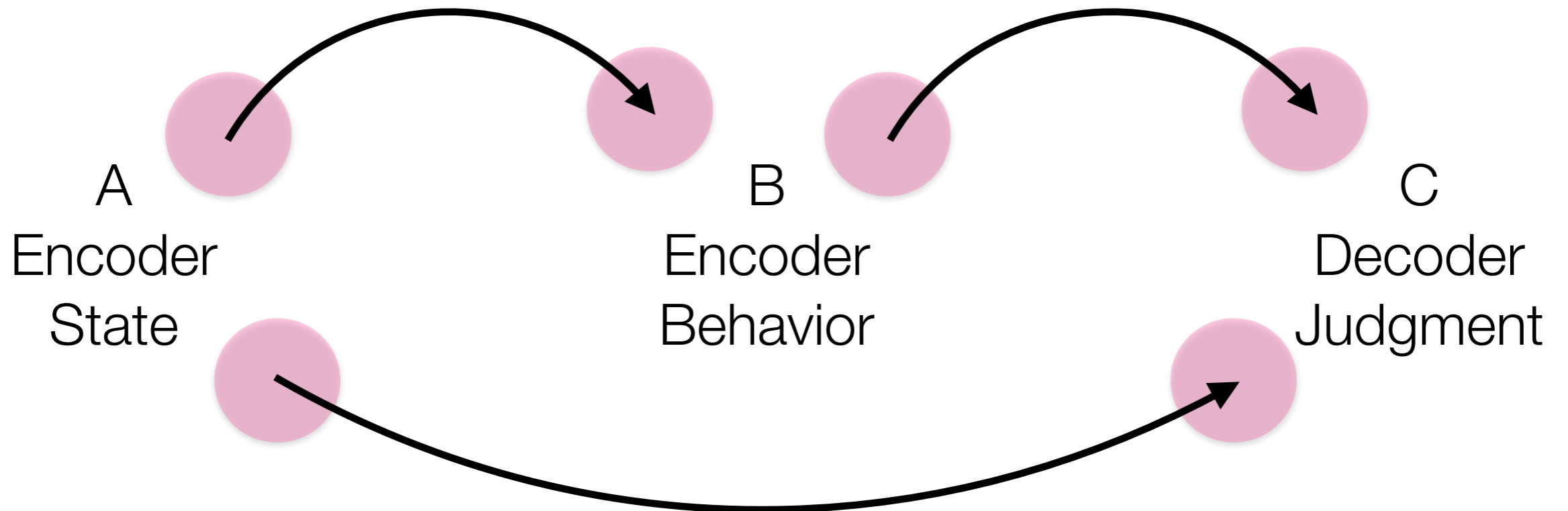
“The term ‘**judgment studies**’ refers most generally to those studies in which behaviors, persons, objects or concepts are **evaluated by one or more judges, raters, coders, or categorizers**, referred to collectively as ‘**judges**’.”

Types of Judgment Studies

Dimensions	Examples
Type of Variable	Dependent vs Independent
Measurement Units	Physical vs Psychological
Reliability	Lower vs Higher
Social Meaning	Lower vs Higher

How does the encoder manifest her/his state through her/his behaviour?

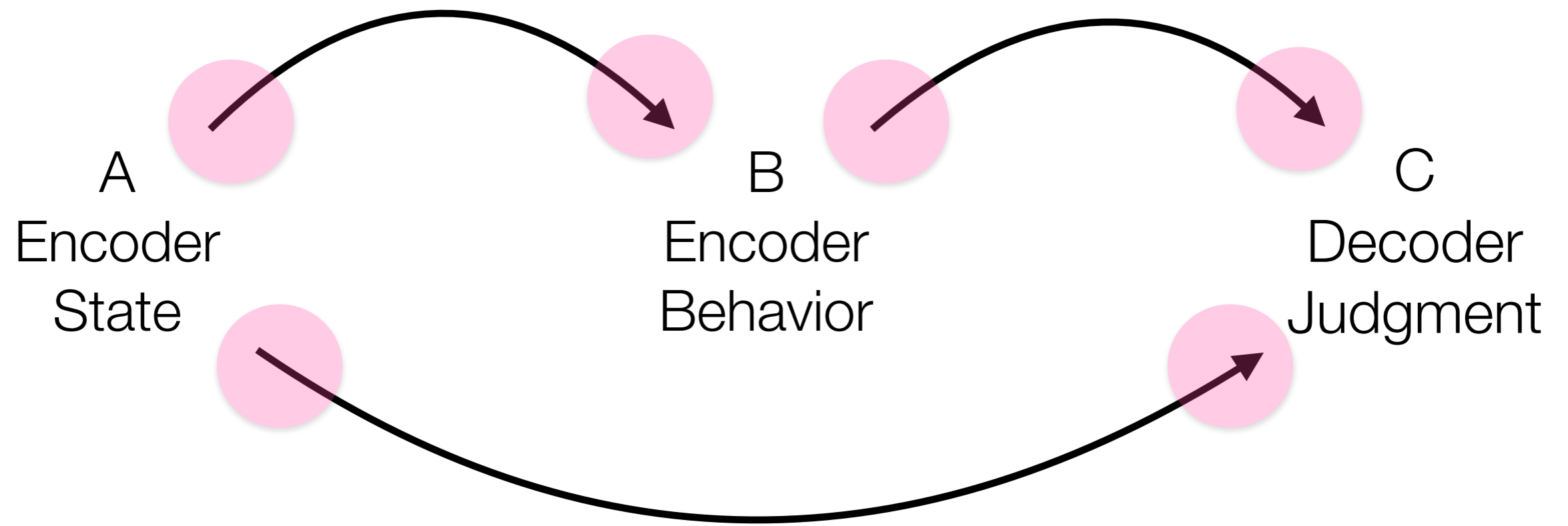
How do decoders detect the behaviour of the encoder?



What is the state the decoder attributes to the encoder?

The State is the independent variable, the behaviour the dependent one

The behaviour is the independent variable, the detection the dependent one



The State is the independent variable, the judgment is the dependent one

Key-Issues in Judgment Studies

- **Selection**: How to select the judges? Should they be expert or naive? Are the demographic characteristics important? Are there risks of self-selection? Do the results depend on the judges?
- **Reliability**: How reliable are the judgments? Are they the result of chance or they reflect a tendency common to all judges? How many judges are necessary to obtain reliable judgments?

The WEIRD Problem

- The acronym **WEIRD** stands for **W**hite, **E**ducated, **I**ndustrialised, **R**ich, **D**emocratic;
- WEIRD people represent **roughly 12%** of the world's population, but they are **over-represented** in psychological studies;
- This bias **affects how we understand human nature**, social interaction, and development, as WEIRD populations have specificities;
- Some phenomena do not generalise to other populations.

The Reliability

- The **reliability** can be thought of as the measure of the **consensus among multiple judges**;
- The **consensus** among multiple judges suggests that there is **consistency between observations and judgments**;
- Judgments are **subjective**, but this does **not** mean that they are **random**;
- **In principle**, the higher the consensus, the higher the reliability.

Average Correlation

$$r = \frac{2 \sum_{i=1}^N \sum_{j=i+1}^N r_{ij}}{N(N-1)}$$

- r is the **average** of the correlations r_{ij} between judge i and judge j ;
- N is the number of judges involved in the study.

Reliability (Spearman Brown)

$$r_{SB} = \frac{Nr}{1 + (N - 1)r}$$

- N is the **number of raters** who assessed the data;
- r_{SB} is the **average correlation** between two raters that assessed the same data;
- r_{SB} is bound between 0 and 1 and **it measures the reliability**, the typical requirement is $r_{SB} \geq 0.7$.

Example

Sample	Judge 1	Judge 2	Judge 3	Total
1	5	6	7	18
2	3	6	4	13
3	3	4	6	13
4	2	2	3	7
5	1	4	4	9

S_{tot}^2

S_2^2

The variance of the scores for one judge

The variance of the total for each encoder

Cronbach's α

$$\alpha = \left(\frac{N}{N-1} \right) \frac{S_{tot}^2 - \sum_{j=1}^N S_j^2}{S_{tot}^2}$$

- N is the **number of raters** who assessed the data;
- S_{tot}^2 is the total variance of the judgments;
- S_j^2 is the variance of judge j .

Pros and Cons

- The **Spearman Brown Coefficient** provides an indication about the correlation between the judgments of two randomly selected Judges;
- However, the Coefficient does not say whether all Judges are correlated to the same extent or not.
- The **Cronbach's α** allows one to avoid calculating a large number of correlations when the number of judges is large;
- The value of the Cronbach's α tends to be similar to the one of the other reliability measures.

Outline

- Judgment Studies
- Example: Personality Perception in AAR
- Conclusions

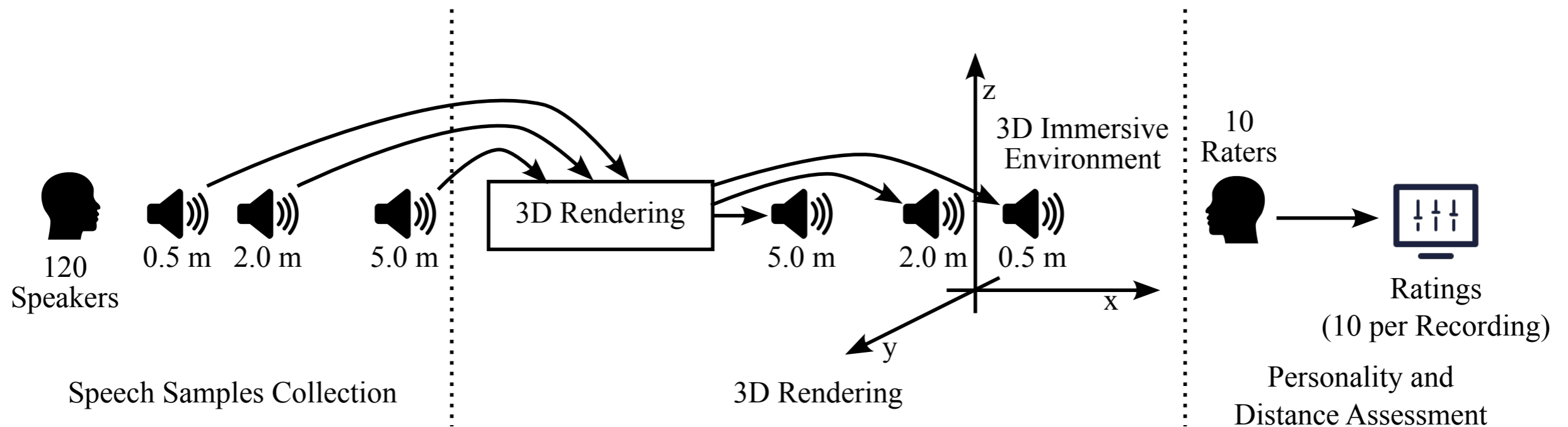
Why Perception?

“People **make social inferences** without intentions, awareness, or effort, i.e., **spontaneously.**”

Why Perception?

“We need to recognise that identification is often most consequential as the categorisation of others, rather than as self-identification.”

The Data



- Every speaker (60 female and 60 male) read **three times the same text** (The North Wind and the Sun), like speaking to **listeners at different distances**;
- The data are rendered in Audio Augmented Reality, an immersive environment that engages only with hearing.

The Data

Gender	Male	Age	Nearfield	Midfield	Farfield	Length
Female	60	31.2±14.0	60	60	60	30 min.
Male	60	33.1±13.5	60	60	60	30 min.
Total	120	32.2±13.7	120	120	120	60 min.

- The recordings are stopped after **10 seconds**;
- The total duration is $10 \times 360 = 3600$ seconds;
- No differences between female and male speakers.

The Data

- The data design **limits the range of factors that change** from one clip to the other;
- The **content is the same** for all clips to avoid variance resulting from the content (what people say);
- The **length is the same** for all clips (10 seconds) to avoid variance resulting from duration;
- The **perceived distance changes** for every clip, the speakers target a specific distance (e.g., 50 cm), but speech production and perception are noisy;
- **Nonverbal speech aspects** (fundamental frequency, loudness, speed, etc.) **change** for every speaker.

Personality

“The **Big Five Personality Factors** appear to provide a set of highly replicable dimensions that parsimoniously and comprehensively **describe most phenotypic individual differences.**”

The Big-Five

- **Extraversion**: tendency to be active, assertive, energetic, outgoing, etc.;
- **Agreeableness**: tendency to be appreciative, forgiving, generous, kind, sympathetic, trusting, etc.;
- **Conscientiousness**: tendency to be efficient, organised, planful, reliable, responsible, thorough, etc.;
- **Neuroticism**: tendency to be anxious, self-pitying, tense, touchy, unstable, worrying, etc.;
- **Openness**: tendency to be artistic, curious, imaginative, insightful, etc.

The Judges

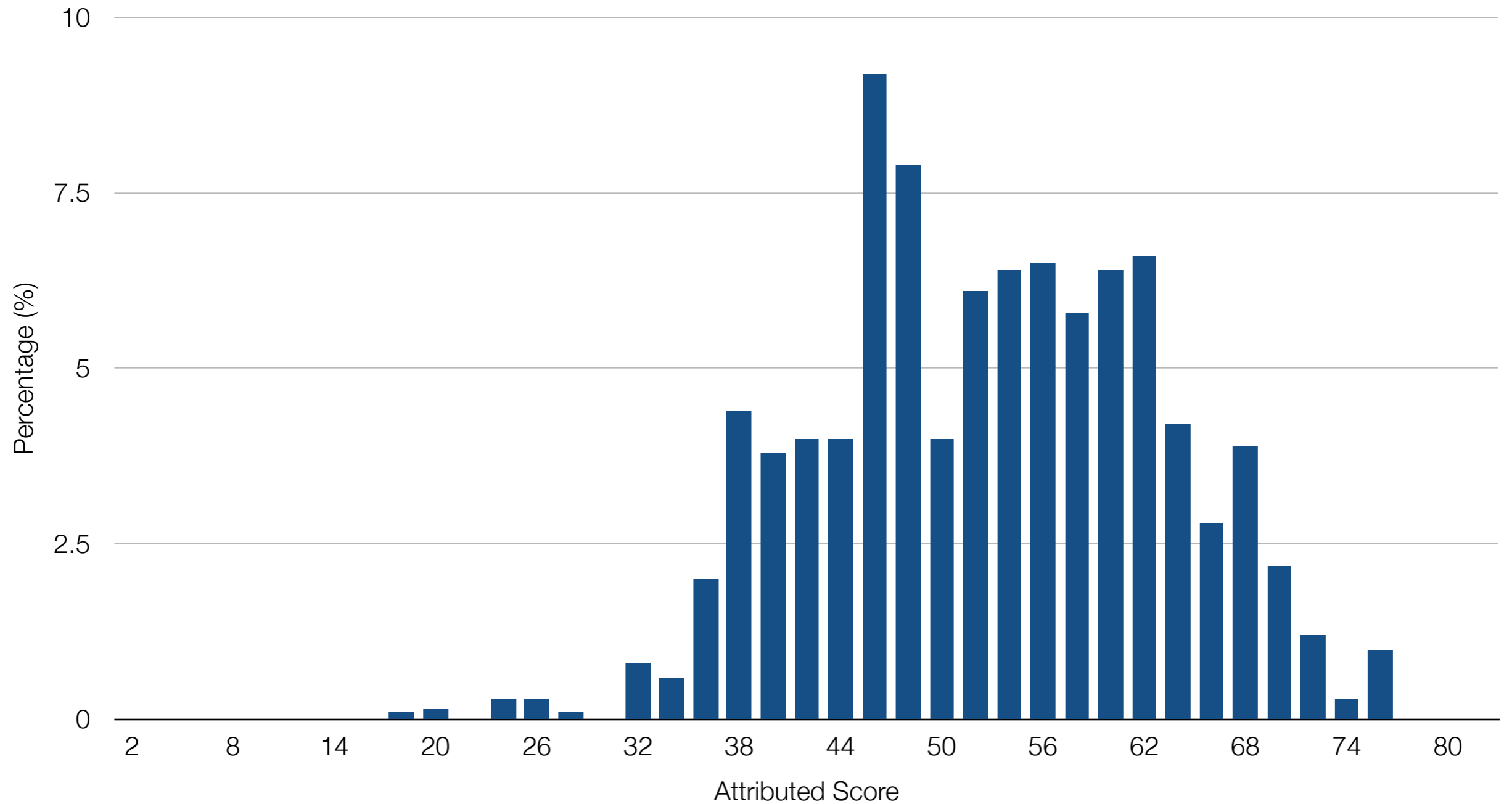
Gender	Male	Age
Female	10	30.0±8.2
Male	8	34.1±12.1
<hr/>		
Total	18	31.8±10.0

- The 360 clips are **randomly split** in non-overlapping blocks of **90 clips** (15 female and 15 male speakers);
- Every **block** is **randomly assigned 10 judges** (5 female and 5 male), some judges assess multiple blocks.

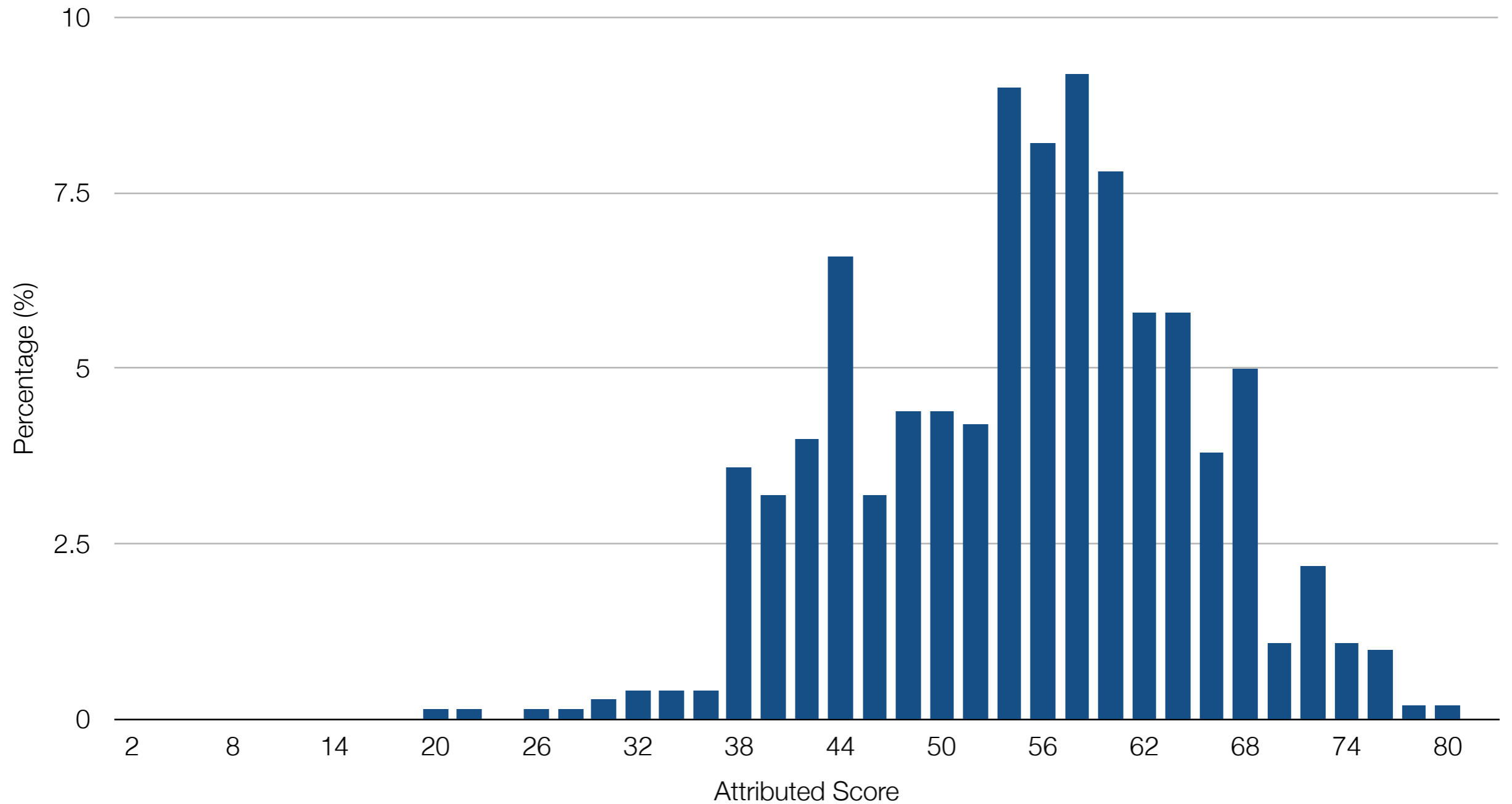
The Task

- The judges were asked to **assess every speaker** in terms of the **Big-Five traits** and the **perceived distance**;
- For the Big-Five, the judges are given a **definition of every trait** and are asked to give a score between 1 and 100;
- The distance must range between **0 and 10 meters**;
- The data were randomly split in four blocks (15 female and 15 male speakers) and every **block** is **randomly assigned 10 judges** (5 female and 5 male), some judges assess multiple blocks.

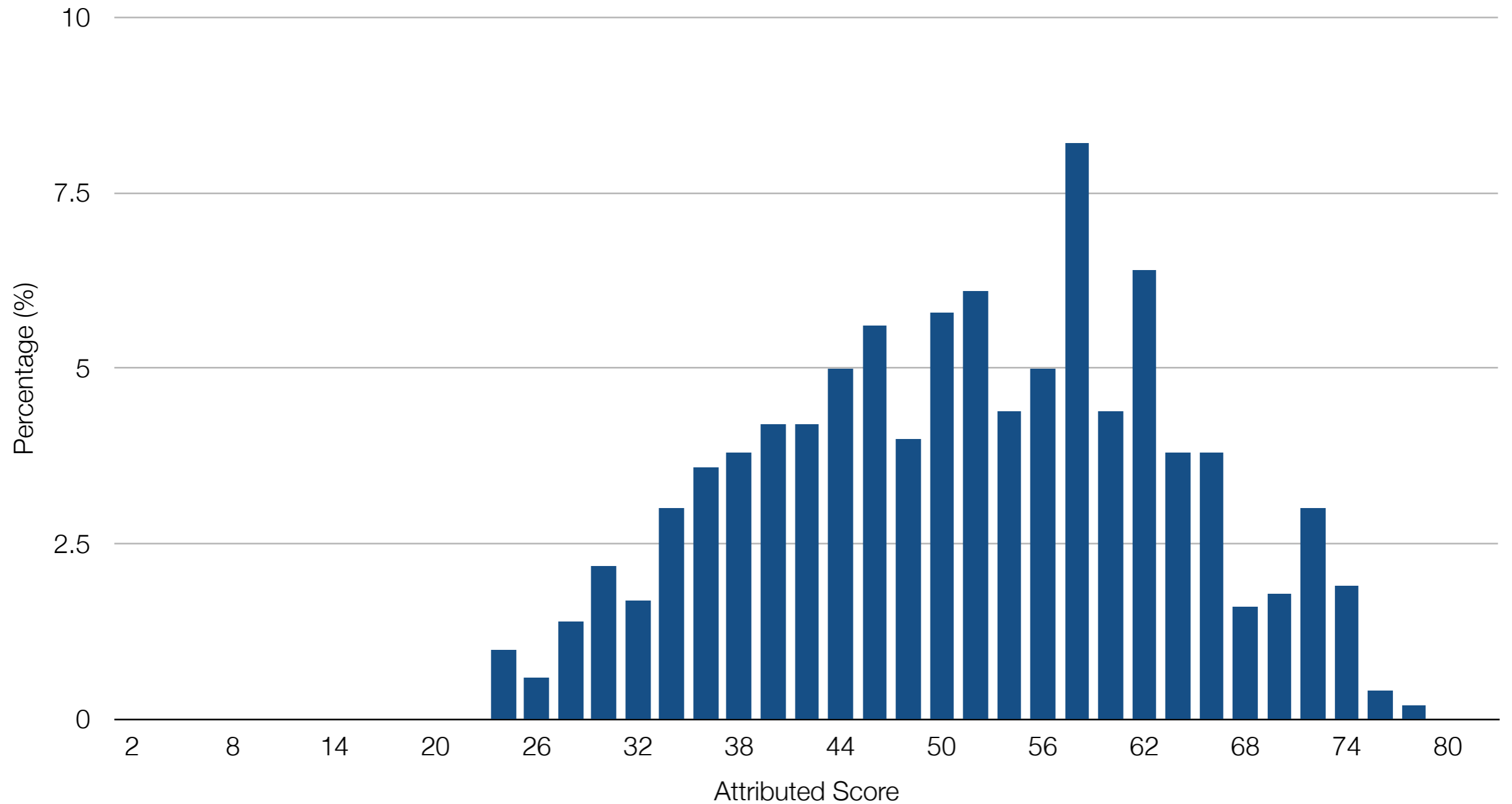
Openness



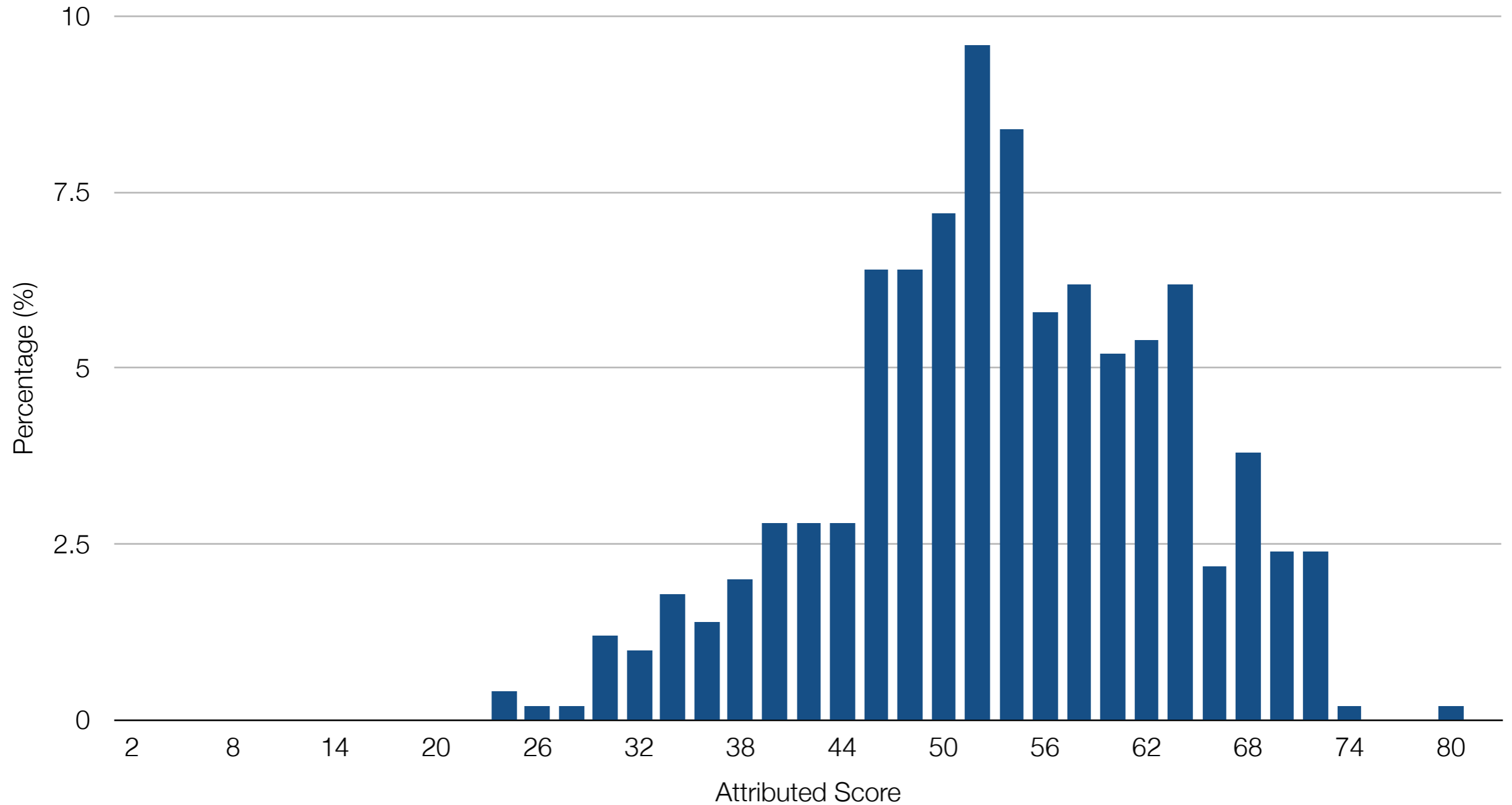
Conscientiousness



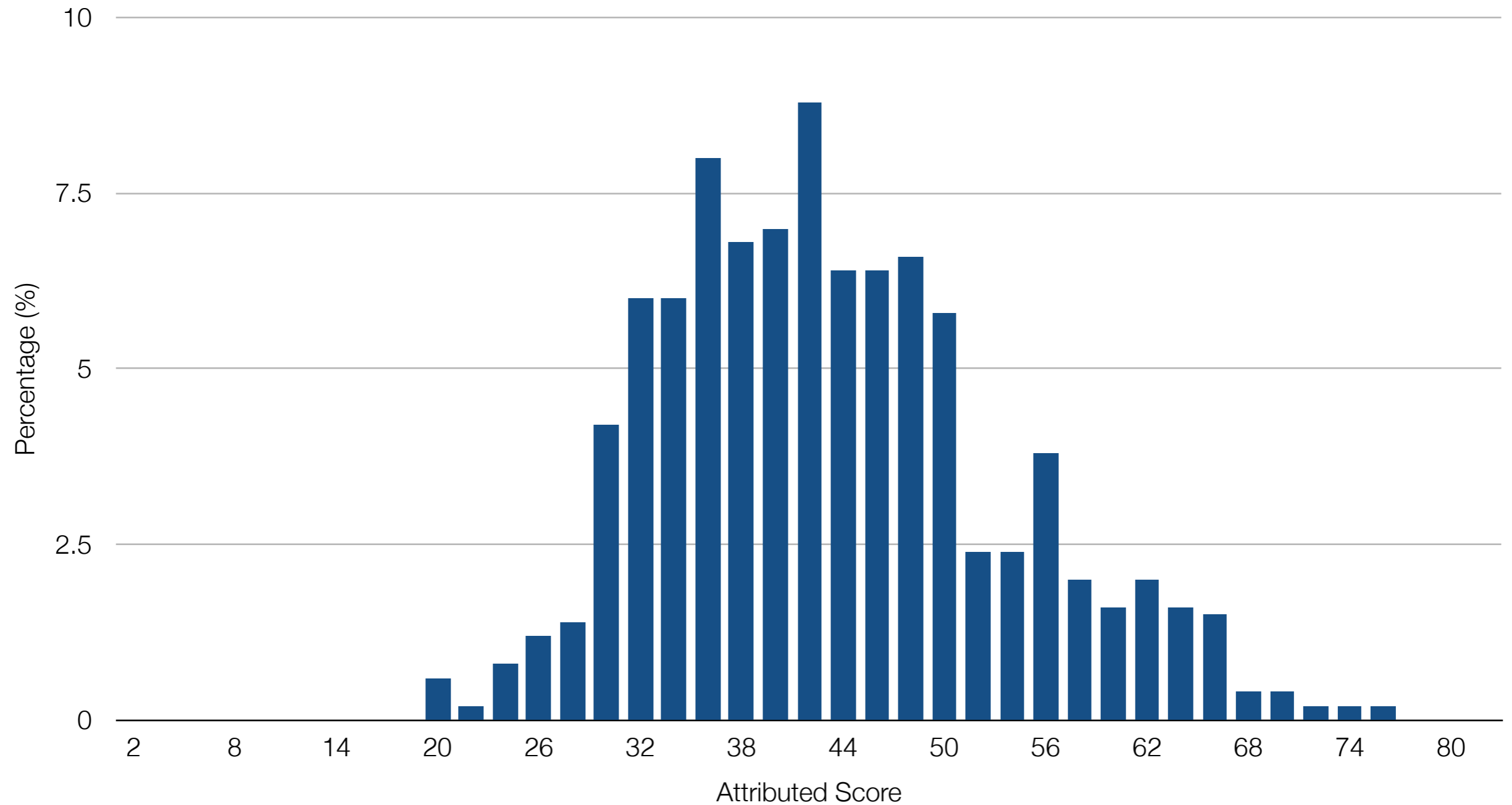
Extraversion



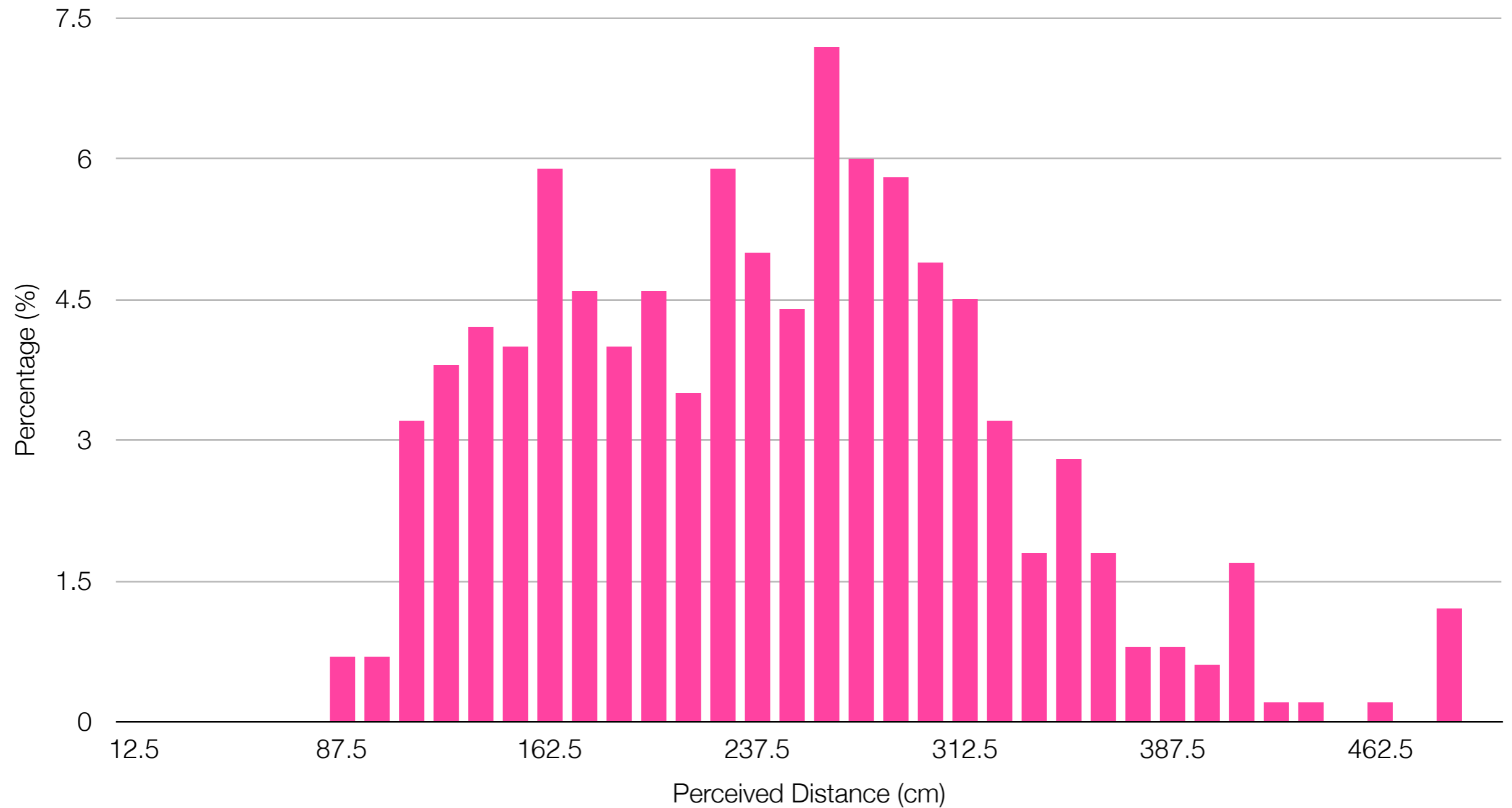
Agreeableness



Neuroticism



Perceived Distance



Reliability

Trait / Distance	Reliability
Openness	0.72
Conscientiousness	0.76
Extraversion	0.84
Agreeableness	0.78
Neuroticism	0.68
Distance	0.87

- The **minimum reliability** for considering the ratings acceptable is conventionally 0.7.

Relative Entropy (I)

$$\sum_{k=1}^N p_k = 1$$

- Consider a **probability distribution** over N mutually exclusive events;
- The value of p_k is the probability of event k taking place.

Relative Entropy (II)

$$H = - \frac{\sum_{k=1}^N p_k \log(p_k)}{\log(N)}$$

- H is the **relative entropy** of the distribution and its value is bound between 0 (one of the p_k is equal to one) and 1 ($p_k = 1/N \forall k \in \{1, \dots, N\}$);
- $H = 0$ means **maximum certainty**, $H = 1$ means **maximum uncertainty**.

Is the Distribution Uniform?

$$H = - \frac{\sum_{k=1}^N p_k \log(p_k)}{\log(N)}$$

- H can be thought of as the **expectation** of the random variable $\log(p)/\log(N)$;
- It is possible to apply a t-test and to verify whether the distributions differ to a statistically significant extent from the uniform distribution.

Relative Entropy

Trait / Distance	Relative Entropy
Openness	0.77
Conscientiousness	0.77
Extraversion	0.80
Agreeableness	0.77
Neuroticism	0.77
Distance	0.73

- According to a t-test, the **entropy is never distinguishable** from the one of the **uniform distribution**.

The Spearman
Correlation
Coefficient

Difference between rank of
trait and rank of GS score
for the same stimulus

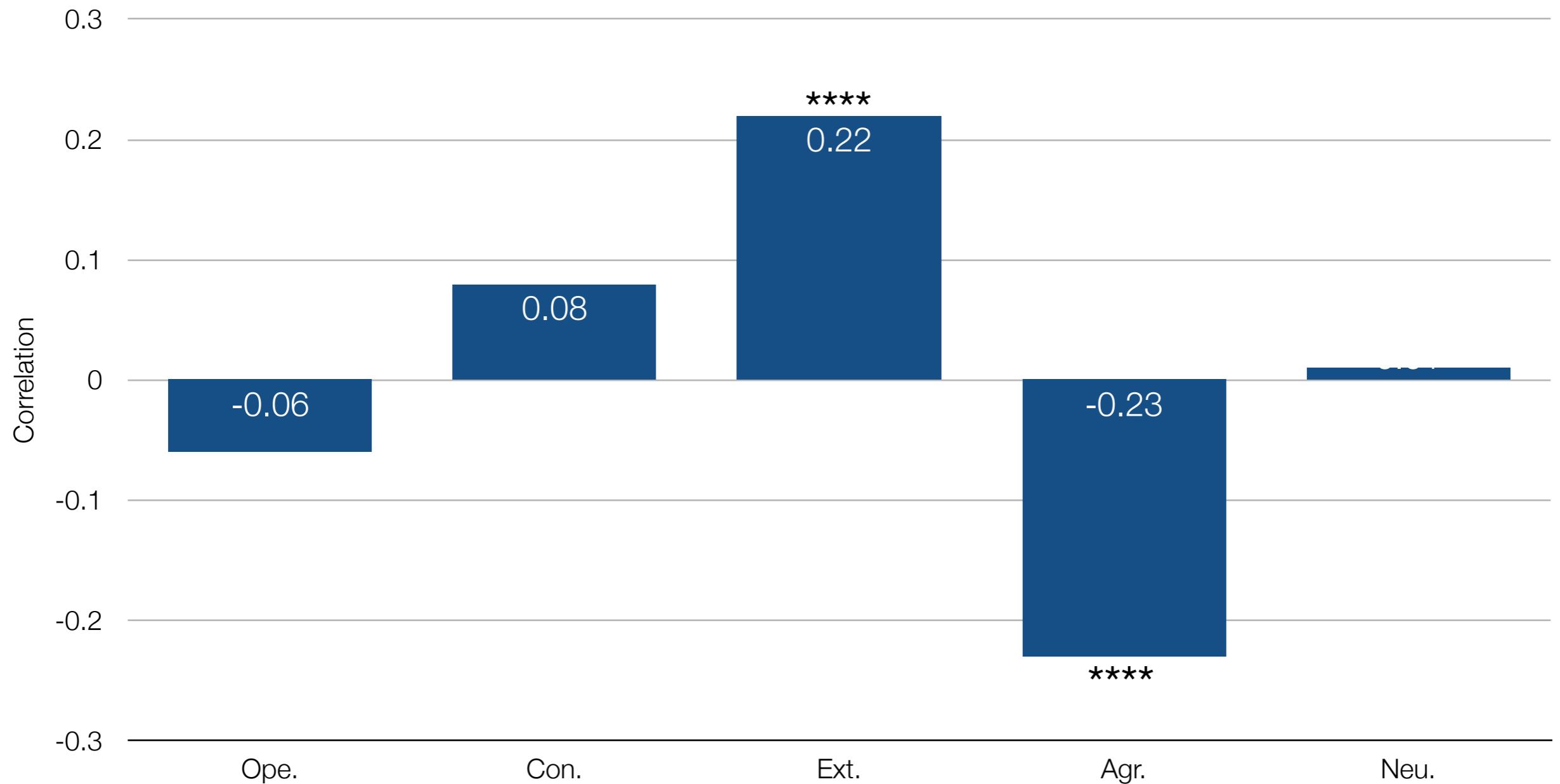
$$r = 1 - \frac{6 \sum_{k=1}^M d(t_k, g_k)}{M(M^2 - 1)}$$

- The Spearman Correlation Coefficient is more robust to outliers than the most common Pearson Correlation.

Correlation Distance-Traits

- The **correlation** between attributed traits and perceived distance can show whether there is a **relationship between the two variables**;
- The main advantage of the correlation is that it is associated to a t variable and, therefore, it is possible to test whether it is **statistically significant**;
- The **Spearman** Correlation Coefficient is **less sensitive to outliers** than the most commonly used Pearson Correlation Coefficient.

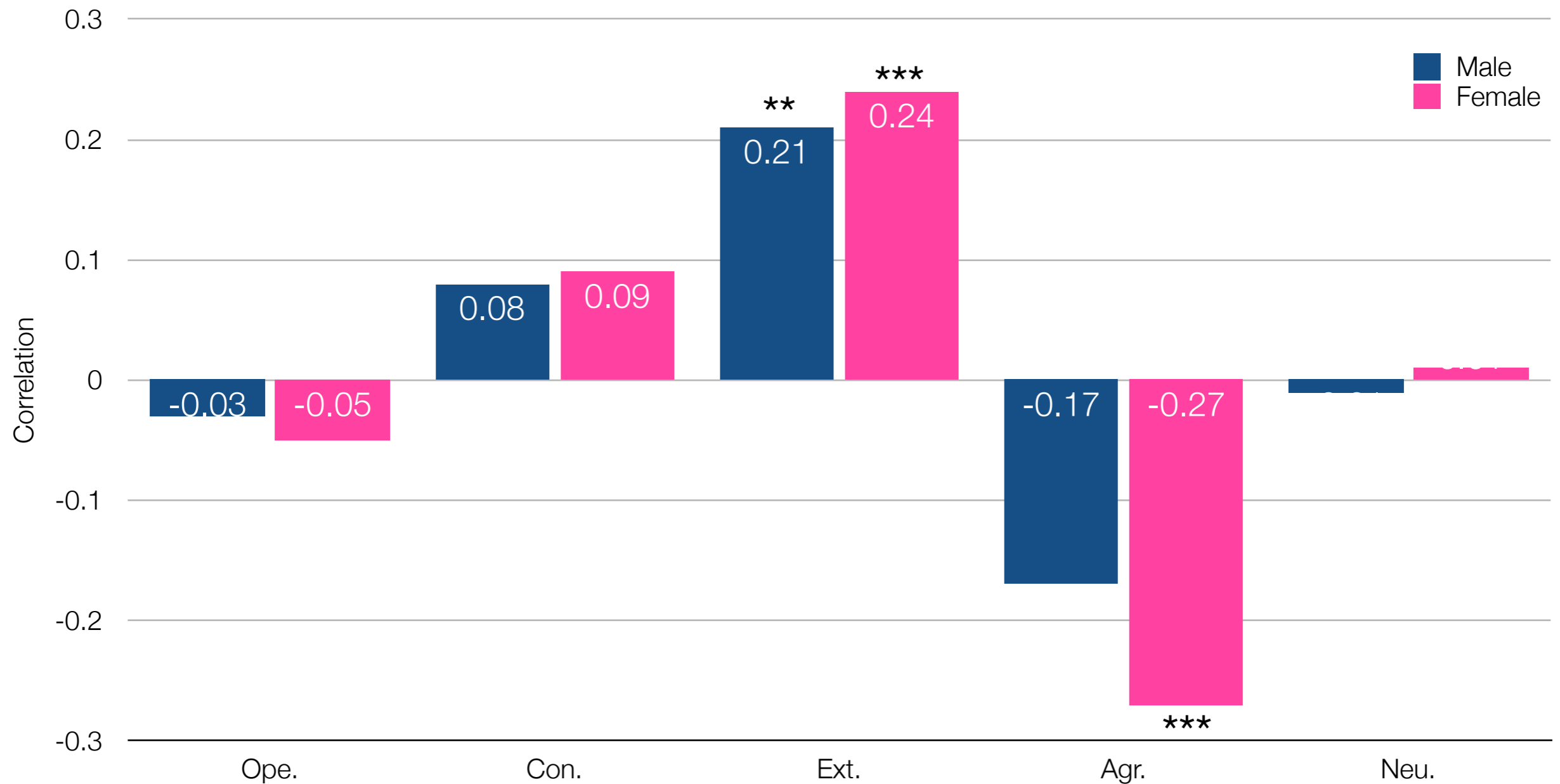
Correlation Distance-Traits



Correlation Distance-Traits

- The correlation is **weak**, but **statistically significant** and **positive** for **Extraversion**: speakers sound extravert when they speak from far away;
- The correlation is **weak**, but **statistically significant** and **negative** for **Agreeableness**: speakers sound agreeable when they speak close;
- Increasing Extraversion reduces Agreeableness, a **tradeoff** must be found **between the two opposite tendencies**.

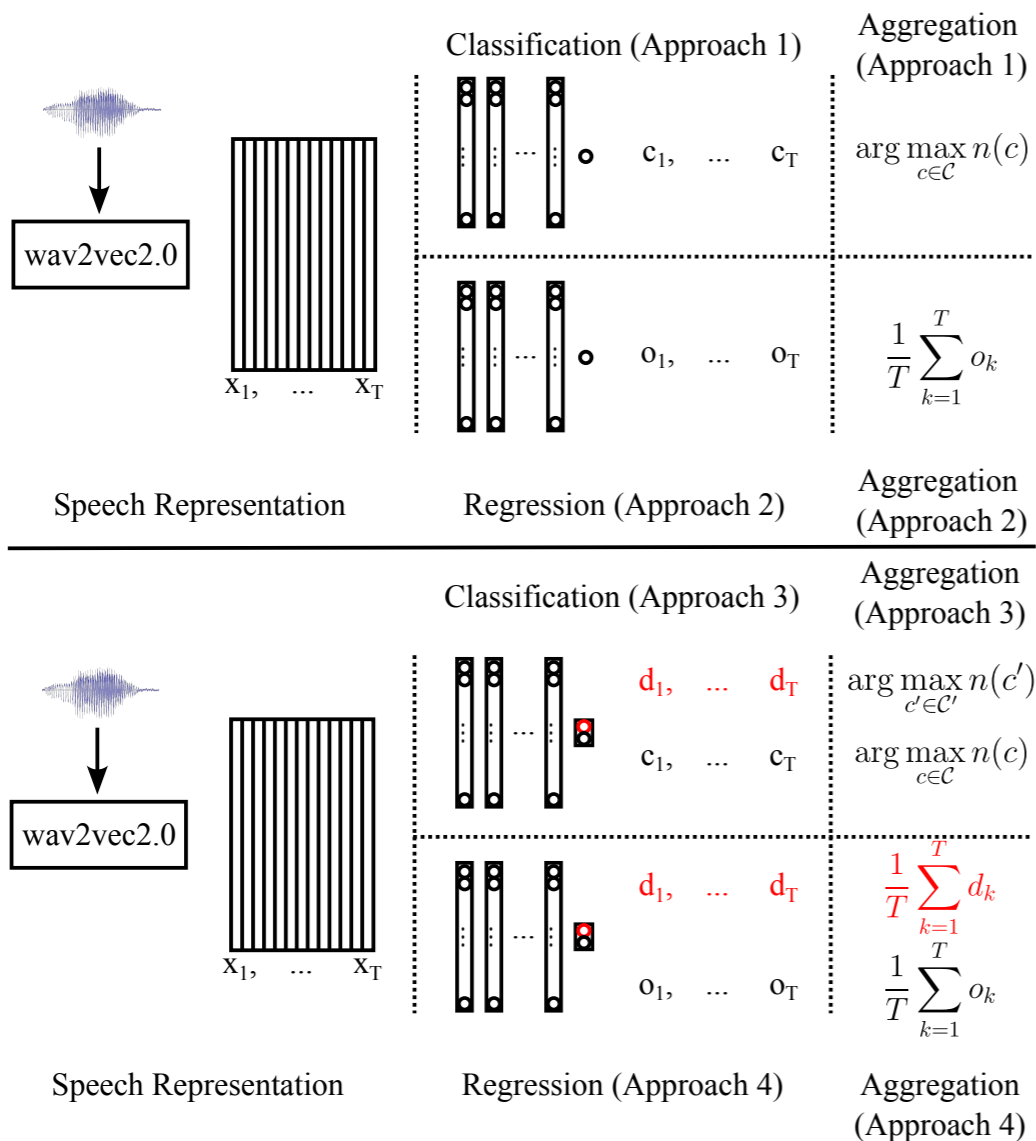
Gender Effects



Gender Effects

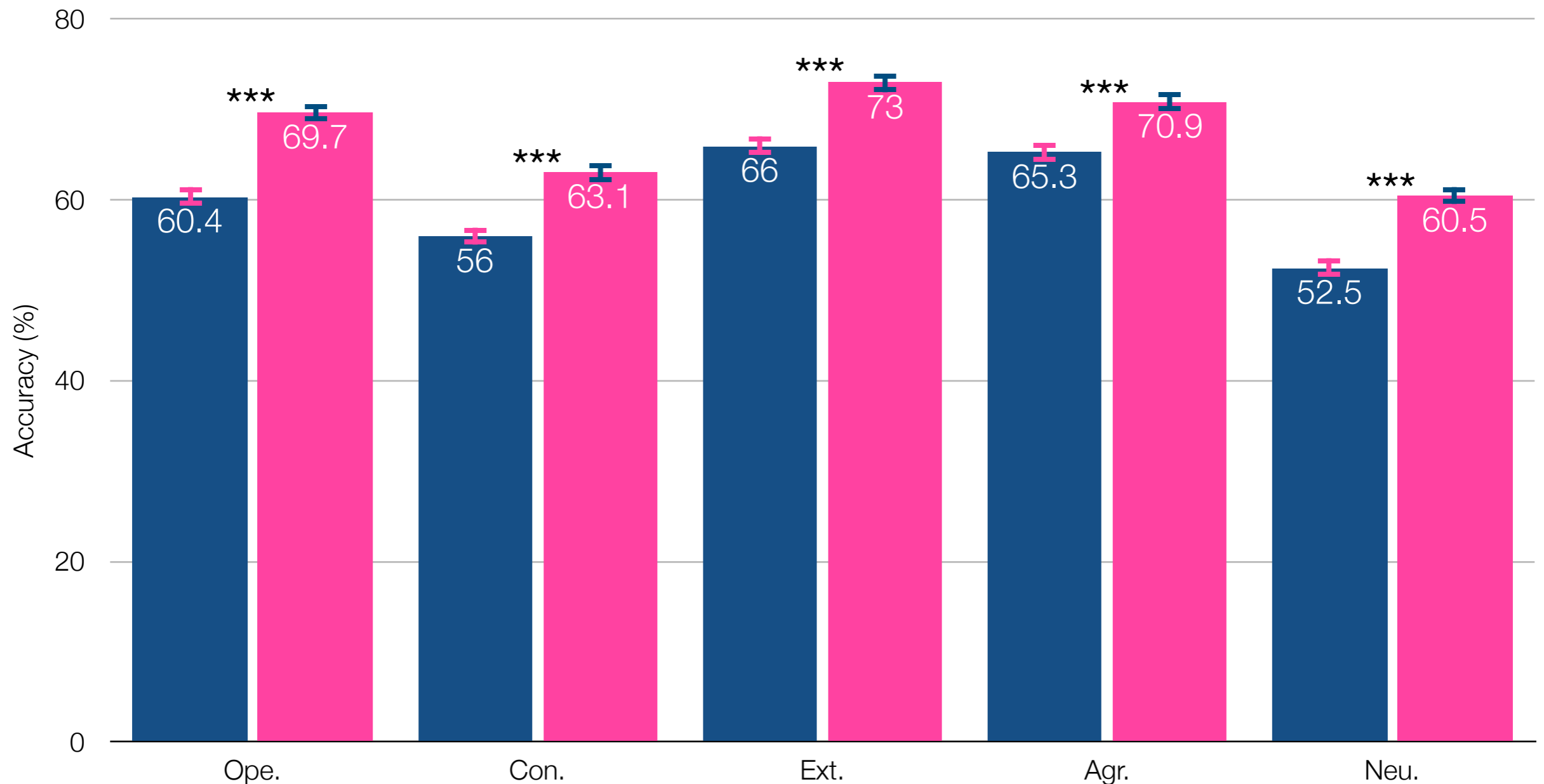
- The correlation is **weak**, but **statistically significant** and **positive** for **Extraversion**, irrespective of gender;
- The correlation is **weak**, but **statistically significant** and **negative** for **Agreeableness**, but the effect applies to female speakers only.
- The problem of the tradeoff seems to apply to female speakers only, but this might depend on the **reduction in the number of speakers**.

Predicting Attributed Traits

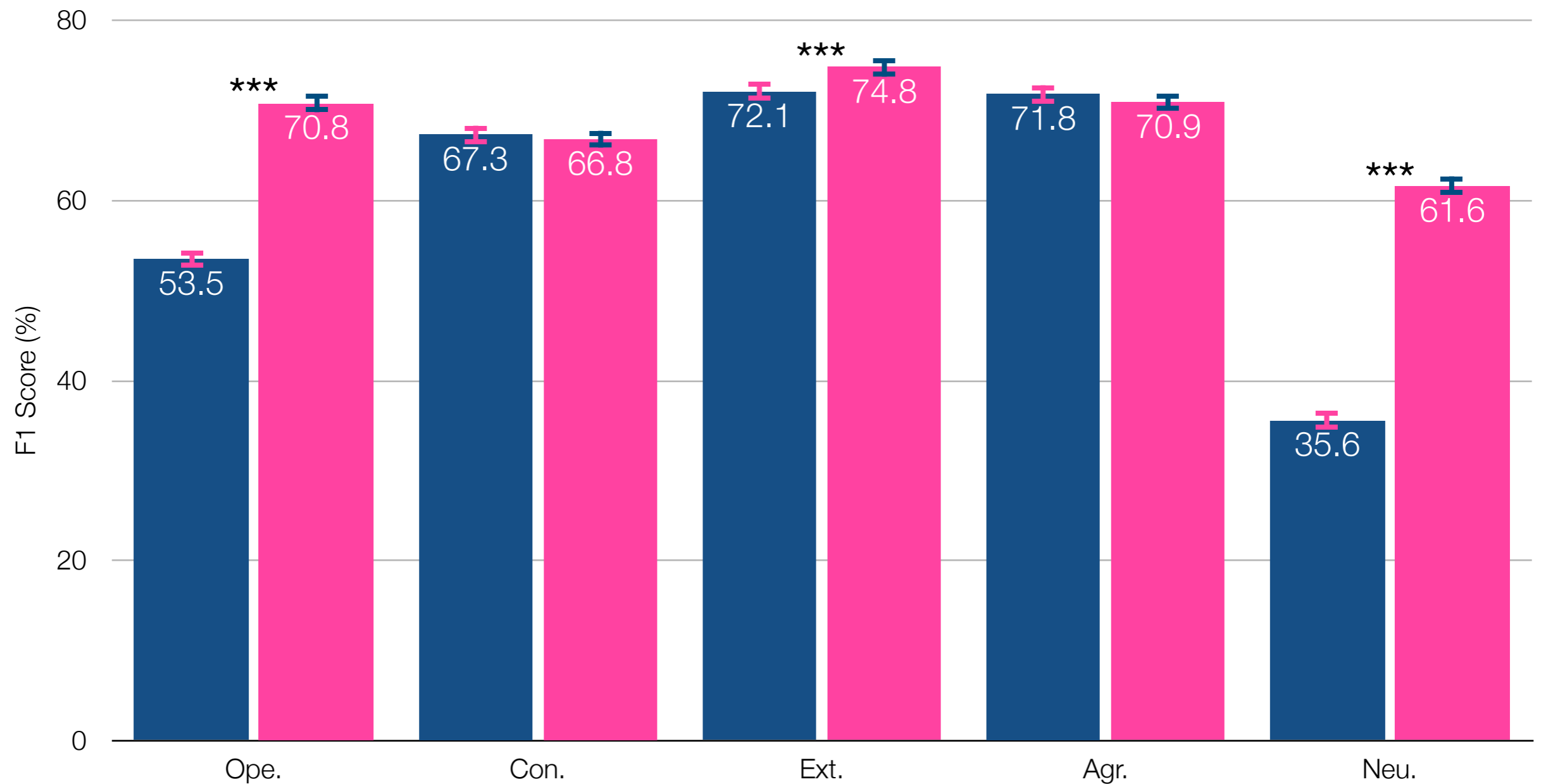


- Predicting whether the speaker is perceived to be above or below median;
- The previous results suggest there is a relationship traits-distance (in some cases);
- Jointly predicting distance and traits should improve the performance.

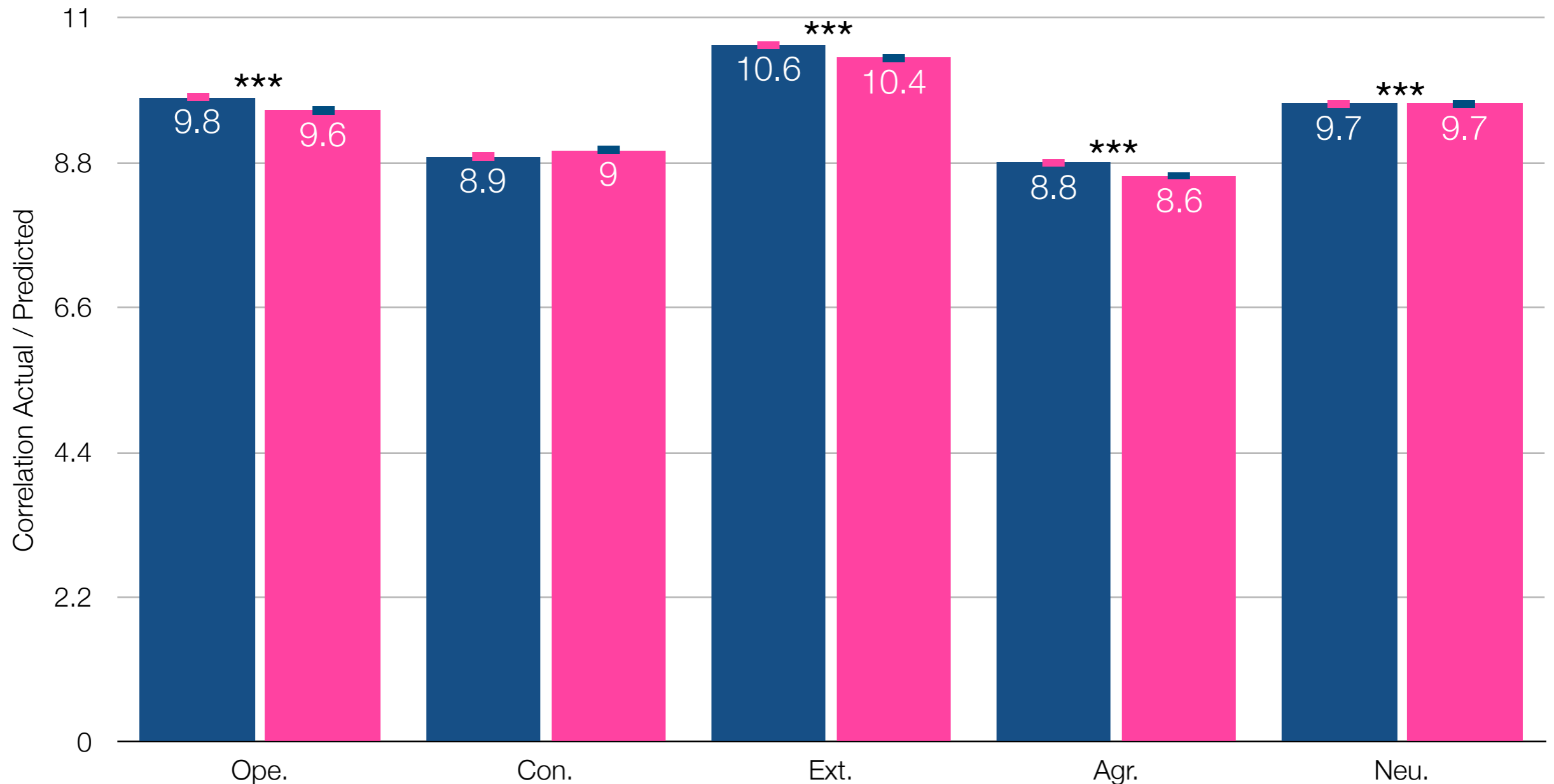
Prediction Results (Accuracy)



Prediction Results (F1 Score)



Prediction Results (Regression)



Prediction

- All traits can be predicted well beyond chance;
- When jointly predicting traits and distance, the accuracy increases for all traits, the F1 Score for Openness, Conscientiousness and Neuroticism;
- The classification improvement seems to confirm that there is a relationship distance-traits, including traits for which the correlation is not significant;
- The improvement suggests relationships distance-traits in cases for which the correlation is not significant;
- An alternative to Psychology?

Outline

- Judgment Studies
- Example: Personality Perception in AAR
- Conclusions

Conclusions

- There is an **interplay** between **perceived distance** and **attribution** of personality **traits**;
- In Audio Augmented Reality, **distance** is one of the **key-perception modalities**;
- **Humans struggle** in assessing the distance of a speaker in absence of visual stimuli and **distance perception is a challenging problem**;
- **Uncertainty** in predicting the perceived **distance** results into **uncertainty in predicting the attributed traits**.

Thank You!