# Acting Data-Driven - But How?

Karsten Lübke, Matthias Gehrke (FOM)

Jörg Horst (FH Bielfeld)

Sebastian Sauer (HS Ansbach)

Gero Szepannek (HS Stralsund)
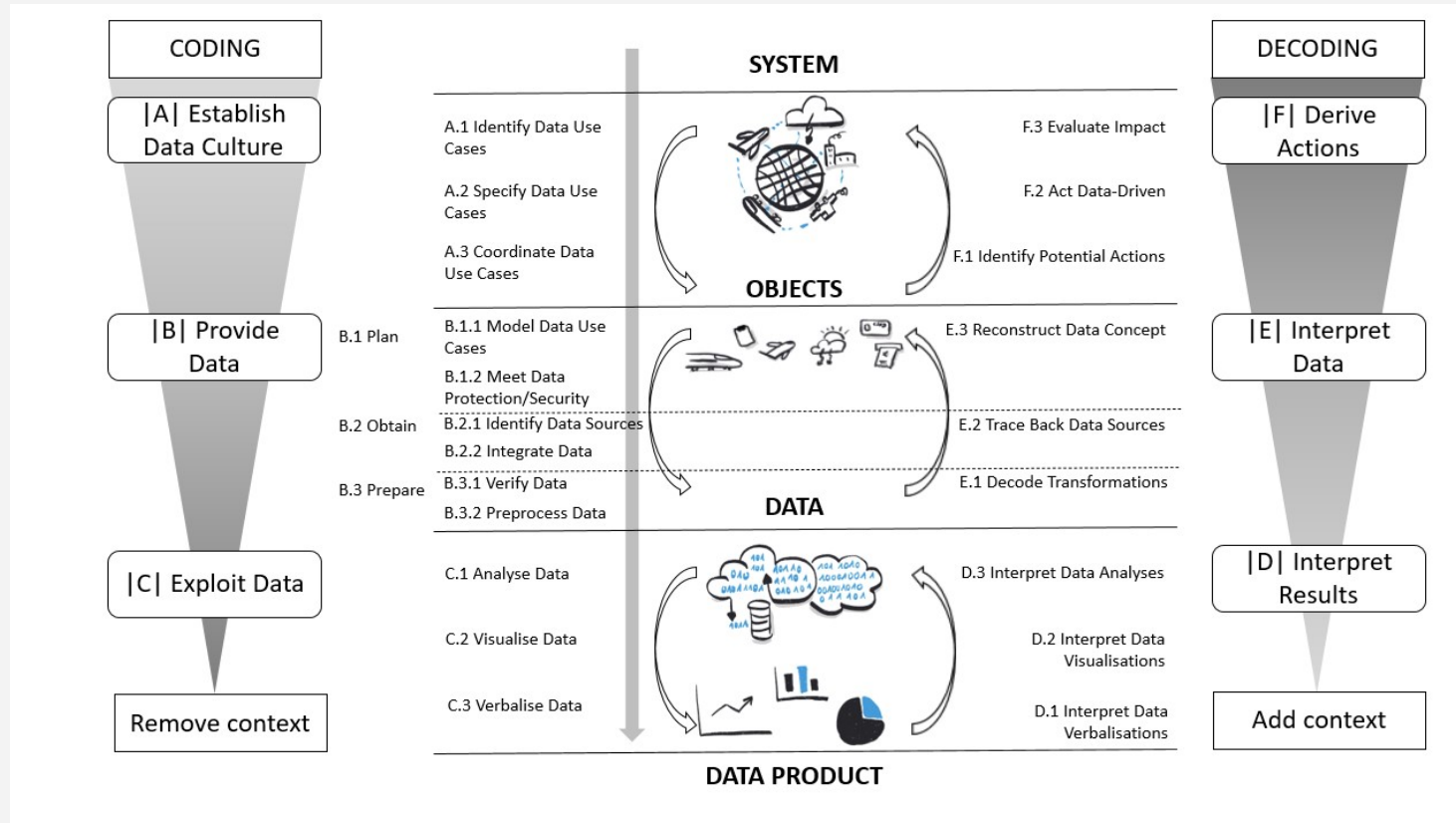
ECDA 2022

# Intro

# Inference from data analysis

Please take part in a (very) short survey: https://bit.ly/30sJNbm

# Data Literacy



Source: Schüller (2020), cf. Data Literacy Charter

# A wobbly bridge

From *A1: Data Use Case* to *F2: Act Data-Driven*:



via GIPHY

# Data science tasks

Hernán et al. (2019) distinguish:

- **Description**: "How can women aged 60–80 years with stroke history be partitioned in classes defined by their characteristics?"

- **Prediction**: "What is the probability of having a stroke next year for women with certain characteristics?"

- **Causal inference**: "Will starting a statin reduce, on average, the risk of stroke in women with certain characteristics?"

# The challenge



Storks Deliver Babies (p = 0.008)

How can we be sure that no human or artificial intelligence does not start colonizing storks to increase birth rate?

# Results

# Back to the Survey: What is inferred?

Structual causal model for data in survey question:
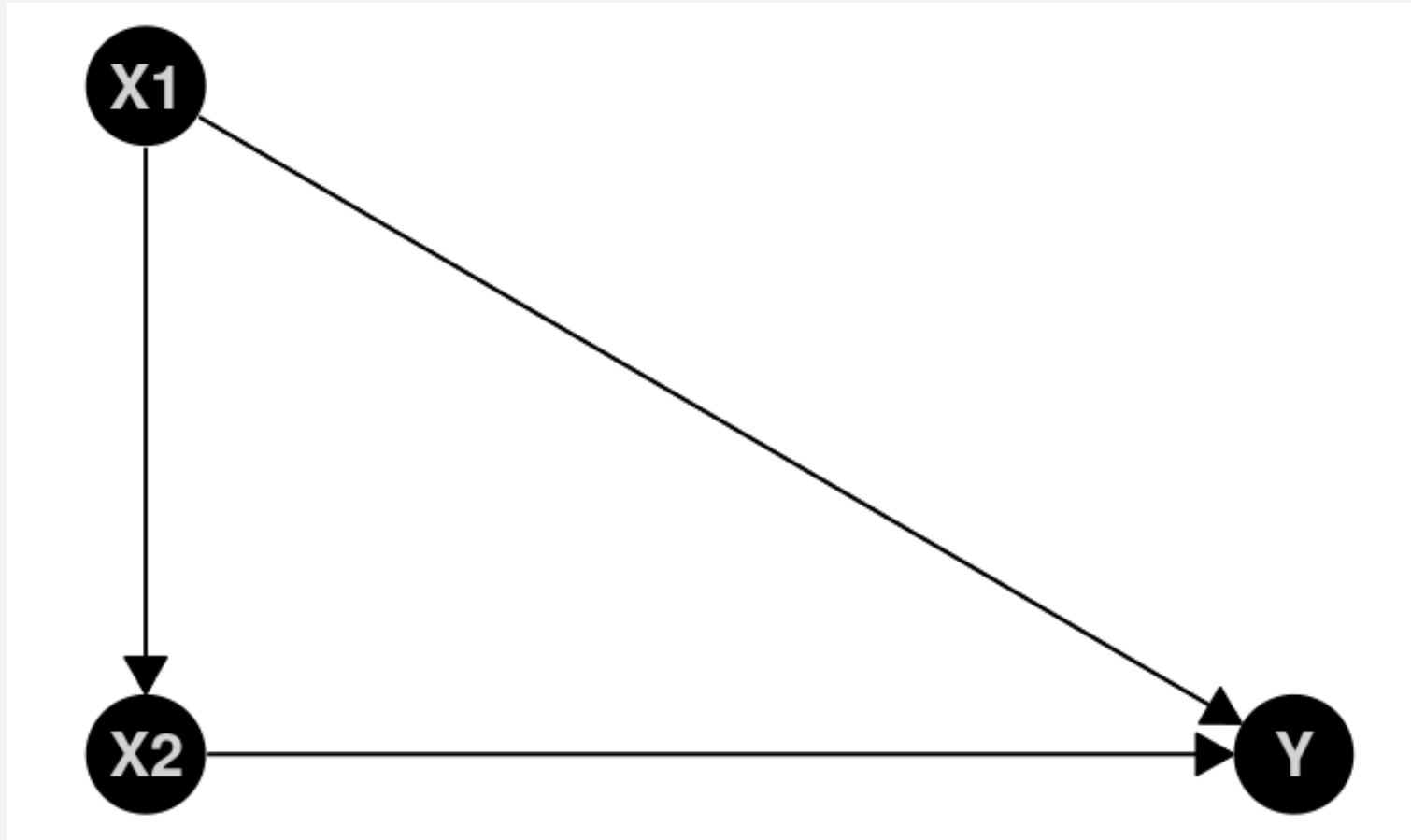
$$X_1 = U_{X_1}, \ U_{X_1} \sim \mathcal{N}(0, \ 10), \quad X_2 = -2X_1 + U_{X_2}, \ U_{X_2} \sim \mathcal{N}(0, \ 1),$$

$$Y = 5X_1 + X_2 + U_Y, \quad U_Y \sim \mathcal{N}(0, \ 5).$$

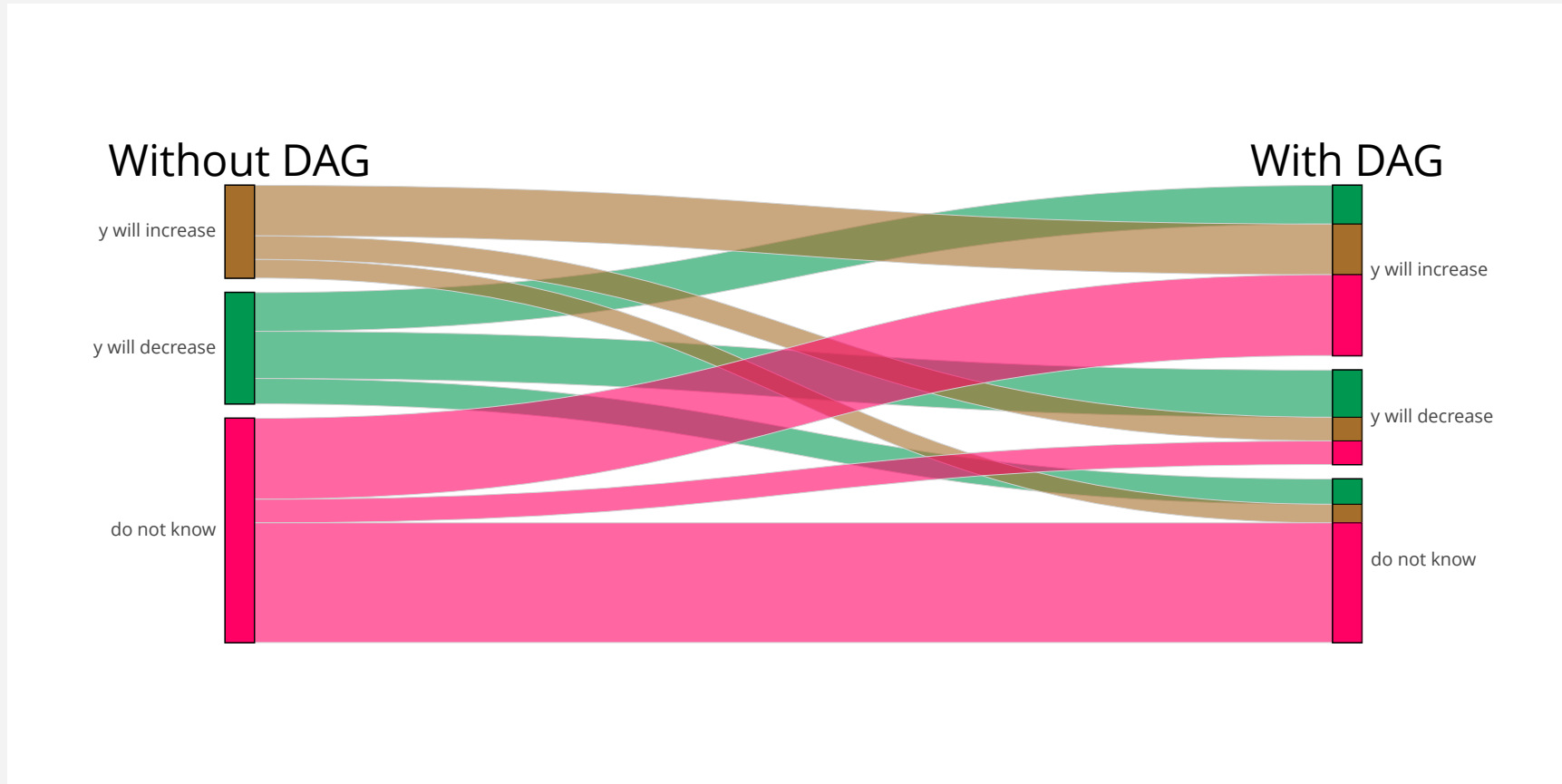Based on linear regression result:

- $\hat{\beta}_2^{(1)} = -1.505$ (excluding $x_1$)

- $\hat{\beta}_2^{(2)} = 0.909$ (including $x_1$)

# DAG

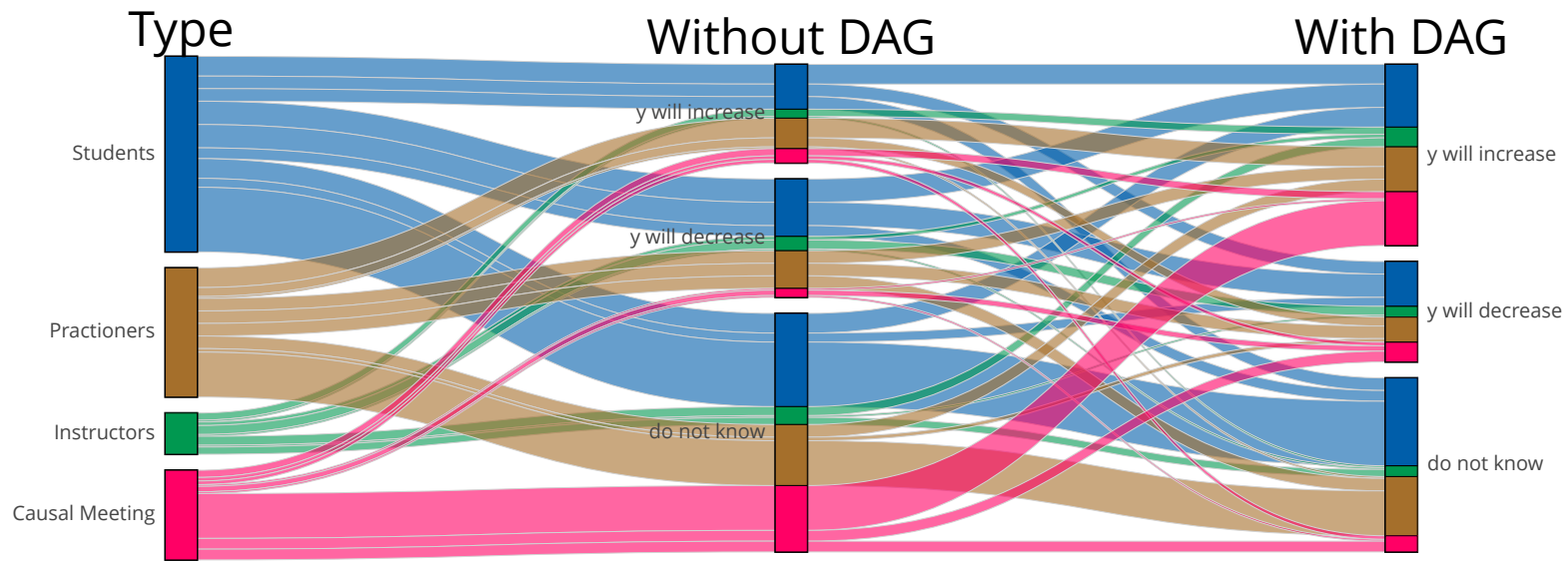# Alluvial diagram

# Alluvial diagram – grouped

# Numerical summary

Correct on both answers:

| Type | n | p.correct |
|---|---|---|
| Causal Meeting | 50 | 0.500 |
| Instructors | 23 | 0.217 |
| Practioners | 72 | 0.097 |
| Students | 109 | 0.101 |

# Freuqentist inference

- For the aggregated data the result is with a p-value of $5.9322569 \times 10^{-5}$ statistically discernible $> 1/9$.

- With a p-value of $5.6886604 \times 10^{-8}$ there are statistically discernible differences between the groups.

# Baysian analysis (uniform prior)

# Outro

# If you've just woken up

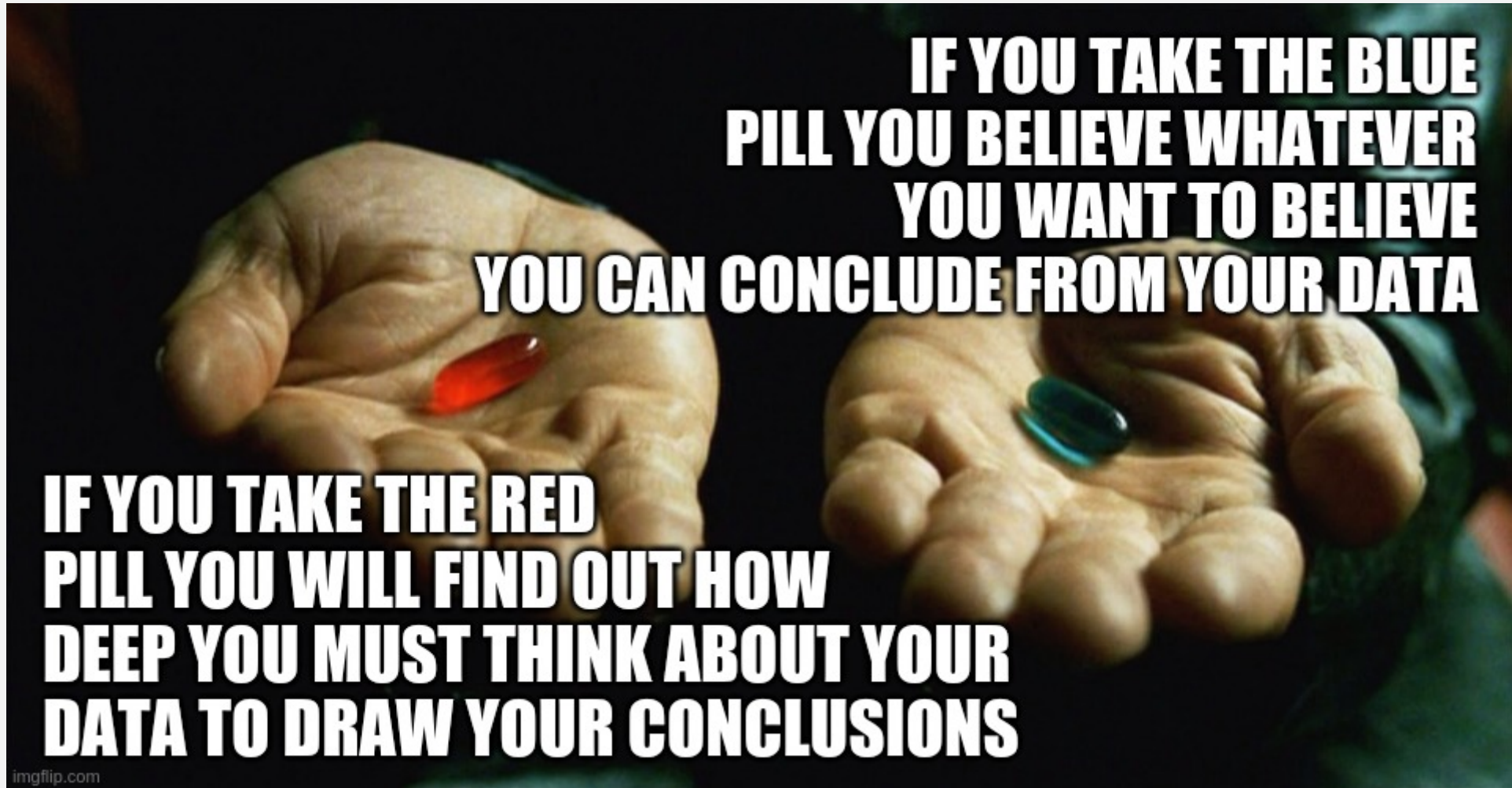Acting Data-Driven - But How?

👉 Far too many draw incorrect conlusions from data analysis. Data analysis skills are not enough to avoid drowning in the data. Integration of DAGs in data science education may be a step in that direction. More research is needed.

# The wrong lesson

Danny Kaplan:

> *What I was saying ...* Data don't speak, they inform our judgment. Interpret data in the context of a whole system.
>
> *What they were hearing ...* The data will say anything you want, depending on how you cut it.

How can we provide a framework to discuss science with data with all stakeholders?

# The End

❤️ Thank you for your participation ❤️

- ✉️: karsten.luebke@fom.de