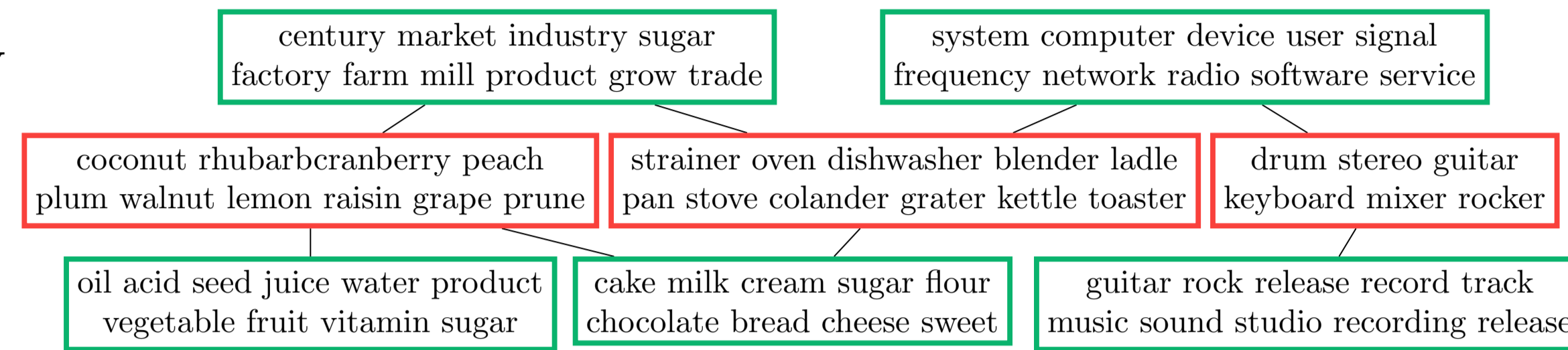# Categorization in the Wild: Category and Feature Learning across Languages

Lea Frermann, Melbourne University, `lea.frermann@unimelb.edu.au`
Mirella Lapata, The University of Edinburgh, `mlap@inf.ed.ac.uk`

## Scaling Models of Categorization I: Categories and Features

- Humans learn **categories** and **features** jointly
- Humans learn **structured** features
- Previous work assumed fixed, relevant features and/or unstructured representations.
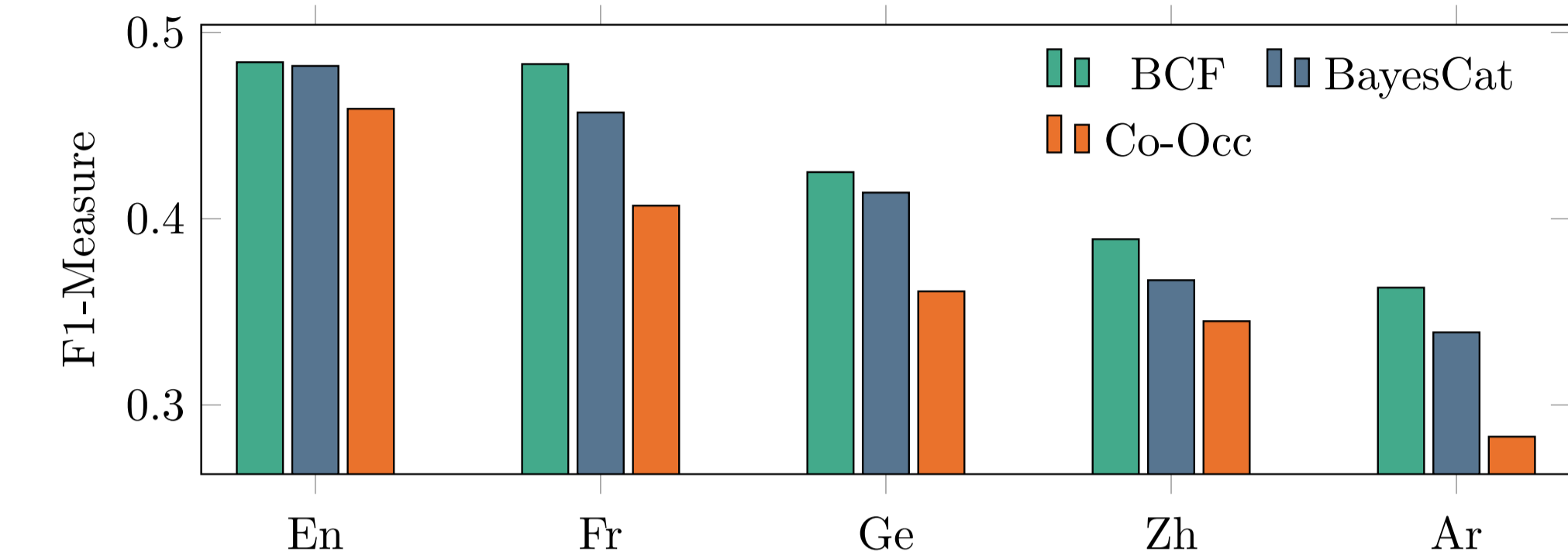
> century market industry sugar factory farm mill product grow trade

> system computer device user signal frequency network radio software service

> coconut rhubarbcranberry peach plum walnut lemon raisin grape prune

> strainer oven dishwasher blender ladle pan stove colander grater kettle toaster

> drum stereo guitar keyboard mixer rocker

> oil acid seed juice water product vegetable fruit vitamin sugar

> cake milk cream sugar flour chocolate bread cheese sweet

> guitar rock release record track music sound studio recording release

## Scaling Models of Categorization II: Languages

- Language as an approximation of the environment
- We apply our models to **five languages**
- **Stimuli**: mentions of *concepts* in linguistic context (*features*)

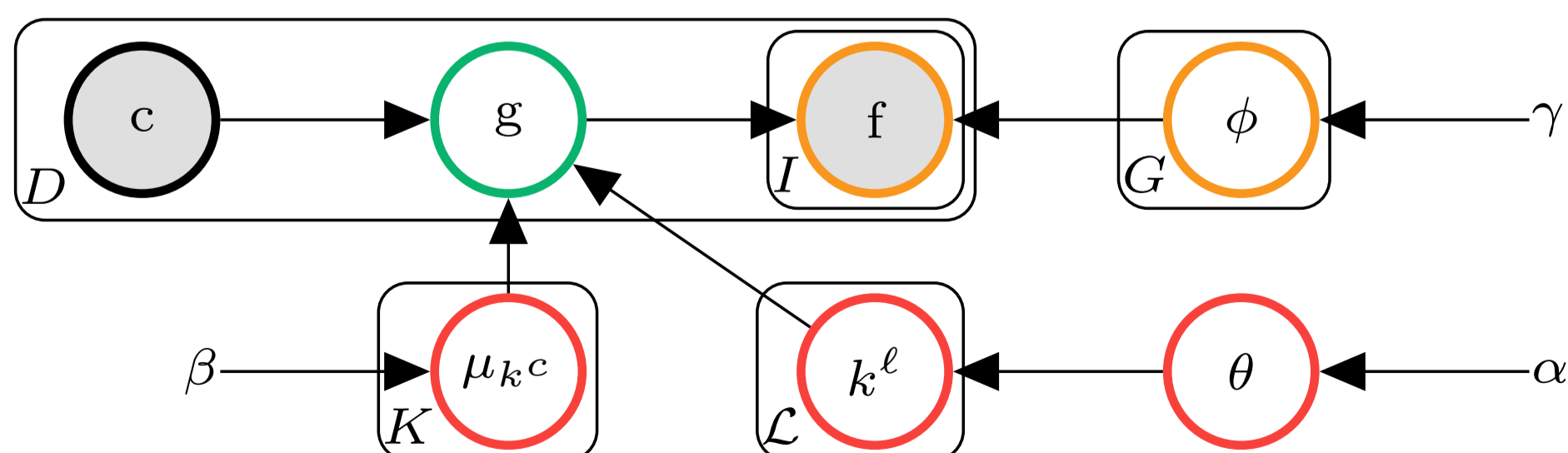| Concept | Natural Language Stimuli | |
|---------|---------------------------|---|
| cat | Les *chats* sont poilus. | 猫有尾巴和爪子。 |
| | *Cats* are carnivores. | Die *Katze* miaut! |
| dog | الكلب لديه الفراء. | Les *chiens* ont des queues. |
| | *Hunde* essen Fleisch. | Look, the *dog* is playing! |
| kiwi | Can you cut me a *kiwi*? | *Kiwis* sind innen grün. |
| | كيويس لديها بذور. | Ce *kiwi* est savoureux. |

## Experiment 1: Category Quality



- **BCF**: our model, **BayesCat**: categorization model with unstructured features, **Co-occ**: co-occurrence model
- Model categorization vs. human-created reference
- Metric: F-1 measure of purity/collocation

## Scaling Models of Categorization III: Diversity

| | En | Fr | Ge | Zh | Ar |
|---|---|---|---|---|---|
| # Concepts | 491 | 484 | 482 | 450 | 394 |
| # Features | 5,898 | 6,416 | 6,981 | 6,516 | 5,870 |
| # Stimuli | 418,755 | 258,499 | 233,175 | 147,386 | 86,908 |

- Hundreds of concepts (from EN feature elicitation studies) Manually translated by native speakers
- In principle unrestricted features (contexts)
- Large sets of stimuli, derived from language-specific Wikipedias
- Concepts, gold categories, stimuli are available `here`

## BCF: A Bayesian Model of Category and Features



- Observe **concept** $c$; Retrieve **category**; Generate **feature type** given category; Generate **features** given feature type
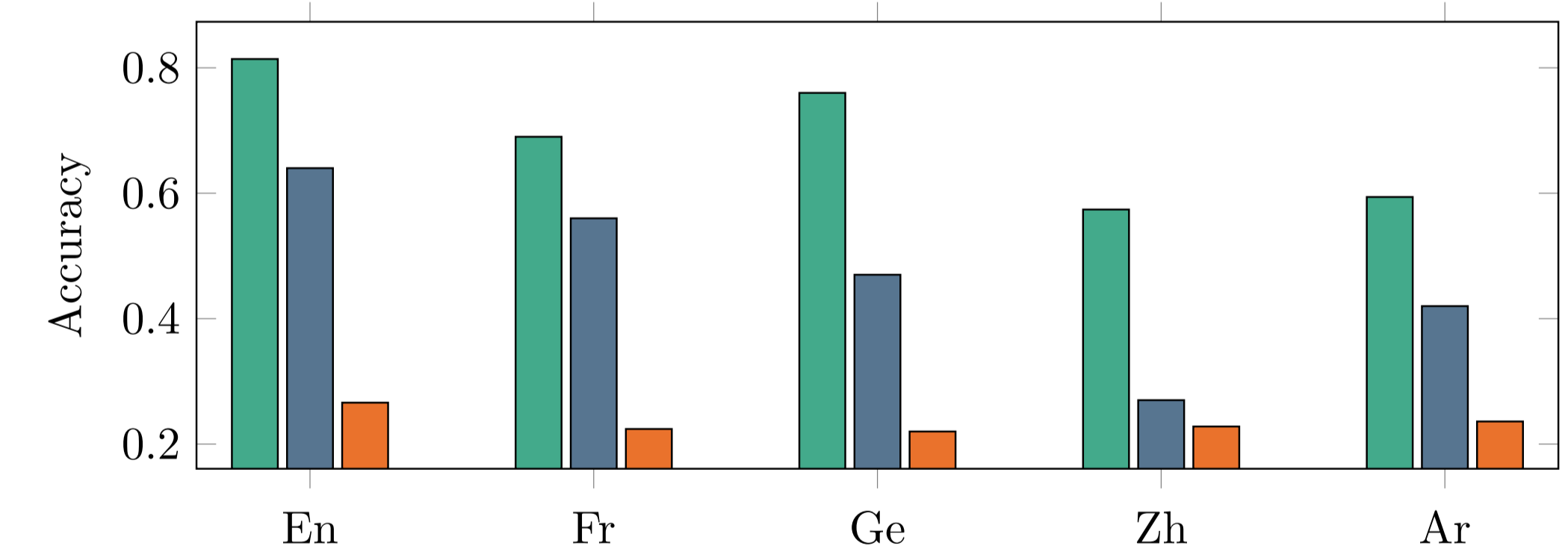- Approximate inference via Gibbs Sampling

## Experiment 2: Feature Quality

### Setup

- Human evaluation through crowd-sourcing; native speakers of the respective languages
- Intrusion paradigm: spot the "intruder" word (feature), which was randomly inserted in the list

### Feature Coherence

| 'Select the intruder word.' | | | | | |
|---|---|---|---|---|---|
| ○ | ○ | ○ | ○ | ● | ○ |
| color | green | blue | white | milk | red |
| ○ | ● | ○ | ○ | ○ | ○ |
| cell | violin | study | protein | human | disease |



### Feature Relevance

| 'Select intruder feature type (right) wrt category (left).' | | |
|---|---|---|
| *wasp ant* | ○ | insect beetle family larva spider |
| *caterpillar* | ○ | tree leaf plant nest grow |
| *hornet moth* | ● | guitar piano clarinet flute |
| *housefly* | | trumpet |
| *beetle* | ○ | male female egg length cm |
| *honeydew* | ○ | white brown dark tail color |
| *grasshopper* | ○ | population habitat bird forest water |