# Data Mining and Landscape Architecting: On Creating Multimodal Resources



#### **Desmond Elliott**

Language and Multimodal Processing Group Department of Computer Science University of Copenhagen

#### The Rise of Culture x Multimodality

• It is important to work with resources that represent the diversity of the cultures for reasons of heritage, inclusion, and equal evaluation



Density map of geographical distribution of images in ImageNet (DeVries+, 2019)



## So how should we collect resources?

#### Two Basic Types of Resource

1. Mined from existing data sources







#### Four Collector Stereotypes

**Data Miner** 

Detectorist

Permaculturist

Baroque









Large-scale extraction

Focused search for artefacts

Minimal intervention

Painstaking design

### **Five Key Factors**

## 1. Scope

What type of data are you aiming to collect?

## 2. Annotator Relationship

How are you working together?

## 3. Images

What type of images?

## 4. Texts

What type of texts?

## 5. Tooling

#### **Use Cases**

#### Multi30K (2015)



- 1. The two men on the scaffolding are helping to build a red brick wall.
- 2. Zwei Mauerer mauern ein Haus zusammen.

#### MaRVL (2021)



(b) *Görsellerden birinde dizlerinde kanun bulunan birden çok insan var.* ("In one of the images, there are multiple people with qanuns on their knees.", concept: *Kanun (çalgı)* (QANUN, a popular instrument in Turkey), label: TRUE)

#### FoodieQA (2024)



#### Kaleidoscope (2025)



- What type of data are you aiming to collect?
- Does your dataset come from a broad collection of concepts / domains, or is it a deep dive into a specific subject matter?



#### Multi30K: Replicate



A brown dog is running after the black dog.

 Ein schwarzer und ein
 brauner Hund rennen auf steinigem Boden aufeinander zu

#### MaRVL: Concepts beyond ImageNet

Concepts







#### FoodieQA: Fine-grained Chinese Food



#### Kaleidoscope: Exams

Geography	Question #5357 Mathematics	Question #9326
'A' চিত্রে প্রদশিত ঘটনাটি কী? রিক্ষিয় ম B	Alice a beaucoup de bl cm x 4 cm. Elle en met dimensions 4 cm x 4 c dépasse. Combien a-t-	locs, tous de dimensions 1 cm x 2 le plus possible dans une boîte de m x 4 cm sans qu'aucun ne eelle rangé de blocs dans la boîte ?
2) অগ্নুৎপাত	1) 6	
3) ক্ষয়	2) 7	
4) ডৰিগল	3) 8	
iocial Sciences   Bengali   Bangladesh	4) 9	
	STEM   French   France	

## 2. Annotator Relationship

- What is the relationship between the collectors and the team?
- Are you working together towards a shared goal, or are they hired?



#### Multi30K: Hired

D#	2016 07:13	
So for me they are very popular as we get here in Austria hardly good task or surveys that it is only an advantage for us that there is the image description 🛞 I do after every 10th task a smoke break otherwise also it gets too steep to me		
And today also again managed	50 of them, which is a great month 😑	

#### Multi30K: Speculations

₽#

2016 07:14

#### 66

I think it's about how different people perceive images ... ... so what is the thing that strikes nearly all or where the focus is. Possibly even something for an AI ... laughing Since many photos which you describe I have for example, the woman at the end of the stairs lying etc.

Yes good chance when I look so what he has been previously employed. I find something fascinating 😂

way, I have just the absolutely is a dark image of a man apparently in a halfpipe and jumps, but you can not identify absolutely with what if skates or board or something else ... I hate something xD

### Kaleidoscope: Collaboration

- Open-science collaboration with an incredible community of early-career scholars
- Co-authorship offered in exchange for collecting data above a threshold



## 3. Images

- What type of images are in your resource?
- From which sources are you collecting them?
- What are the licenses of the images?
- How are you protecting PII?



#### FoodieQA: Private Images





## 4. Texts

- What type of text are you collecting?
- Where are you getting the text from?
- Which languages are you covering and why?





#### Kaleidoscope: Natural Questions

 The texts were created by exam-board professionals for educational purposes



Consider the Deterministic Finite-state Automaton (DFA) A shown below. The DFA runs on the alphabet {0, 1}, and has the set of states {s, p, q, r}, with s being the start state and p being the only final state. Which one of the following regular expressions correctly describes the language accepted by A?



#### FoodieQA: Procedural Single-Image QA

• The fine-grained taxonomy and careful labelled meant that we could automatically create questions using templates

<dish>是哪个地区的特色菜? (What region is <dish> a specialty dish of?)
<dish>是哪个地区的特色美食? (In which region that <dish> is a local specialty?)
去哪个地方游玩时应该品尝当地的特色美食<dish>? Which place should you visit to taste the local
specialty food <dish>?

以下菜品是哪个地区的特色菜? Which region is this food a specialty?



(Jiangsu)
 京津 (Beijing & Tianjin)
 香港 (Hong Kong)

## 5. Tooling

- Which platforms are you using to create your resource?
- Are you relying on re-usable tools or is your pipeline custom-built?

Re-usable Custom

#### Acknowledgements





- E. Bugliarello R. Ramos





L. Specia



L. Barrault





N. Collier



G. Chrupała



C. Qui



D. Oneață



S. Frank



R. Sanabria





F. Liu



**B.** Martins



E. Hasler



E. Ponti



S. Hooker





A. Alishahi



I. Salazar



### Conclusions

- Two basic types of data in Culture x Multimodality
  - Data that is **mined** from something that already exists
  - Data that is **created** from scratch given a few ingredients and tools
- Lots of factors to keep-in-mind when we are creating a high-quality resource that will be re-used by hundreds-thousands of people
- Remember: not all data should be extracted (Bird, 2021)

#### References

- I. Salazar\*, M. Fernández Burda\*, S. Bin Islam\*, A. Soltani Moakhar\*, S. Singh\*, F. Farestam\*, A. Romanou\*, et al. Kaleidoscope: In-language Exams for Massively Multilingual Vision Evaluation. arXiv:2504.07072.
- W. Li, C. Zhang, J. Li, Q. Peng, R. Tang, L. Zhou, W. Zhang, G. Hu, Y. Yuan, A. Søgaard, D. Hershcovich, and D. Elliott. FoodieQA: A Multimodal Dataset for Fine-Grained Understanding of Chinese Food Culture. EMNLP '24.
- F. Liu\*, E. Bugliarello\*, E. M. Ponti, S. Reddy, N. Collier, and D. Elliott. Visually Grounded Reasoning across Languages and Cultures. EMNLP '21.
- D. Elliott, S. Frank, K. Sima'an, and L. Specia. **Multi30K: Multilingual English-German Image Descriptions**. Workshop on Vision and Language at ACL 2016.