# Upshot from Dayton: Found in Translation

Ranking Arabic Documents with Machine Learning

- 28 Letter Alphabet
- Logical structure (nouns and verbs) similar to English
- Written and read from right to left.
- Written in cursive, making text recognition challenging.

# *Concept of Operation*

- Understand the value of the document
- Limited Interpreters - upshoot documents that have urgent relevance
- Mobile app translation is nice to have

# *Problem Statement*

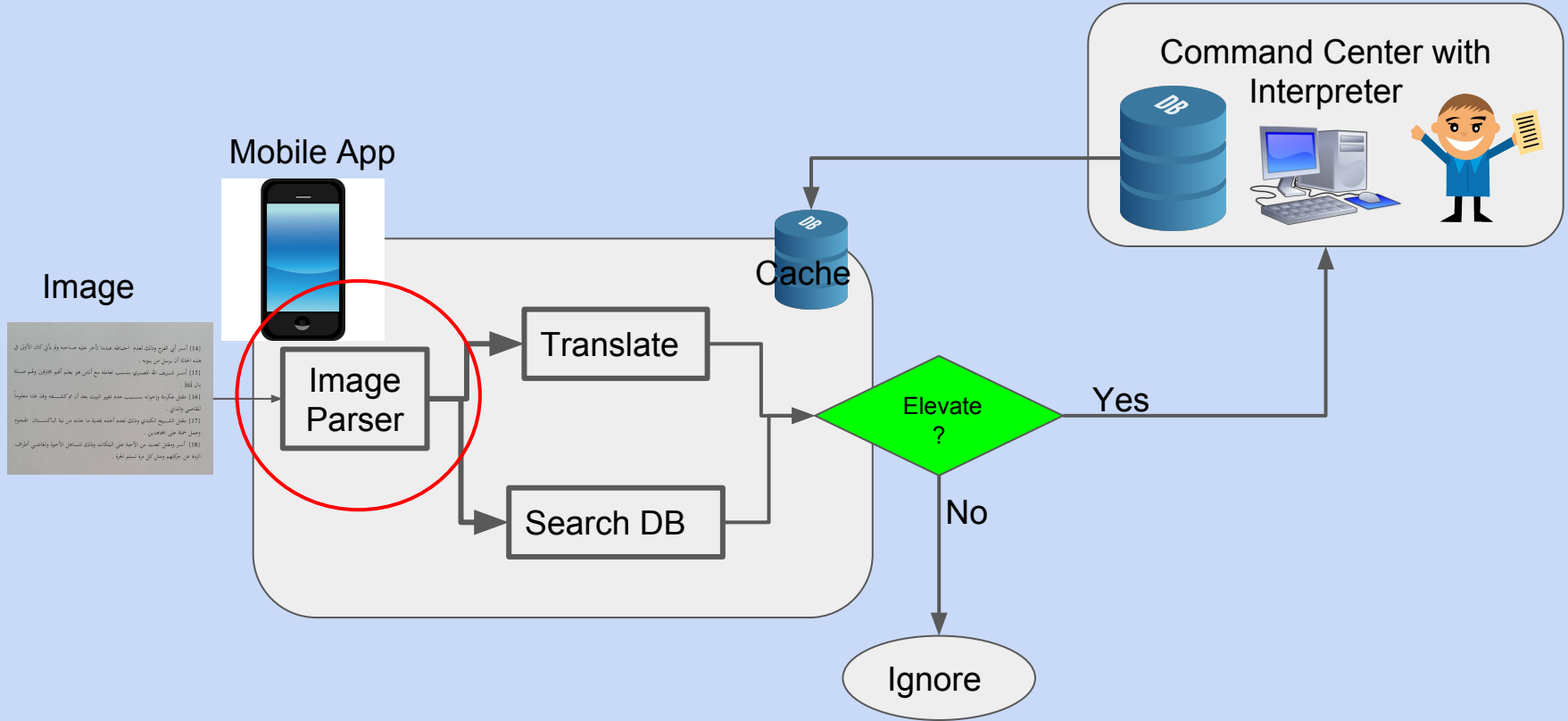- Given a document photo, what other documents is it like?
- Should someone translate it?

١٠/١٣ pt صدى     ١٠/١٣ pt Seria Italic

من عداد الصنائع الإنسانية وهو رسوم وأشكال حرفية
تدل على الكلمات المسموعة الدالة على ما في النفس.
فهو ثاني رتبة عن الدلالة اللغوية وهو صناعة شريفة إذ
الكتابة من خواص الإنسان التي يميز بها عن الحيوان.
وأيضاً فهي تطلع على ما في الضمائر وتتأذى بها
الأغراض إلى البلد البعيد فتقضي الحاجات وقد دفعت
مؤونة المباشرة لها ويطلع بها على العلوم والمعارف
وصحف الأولين وما كتبوه في علومهم وأخبارهم فهي
شريفة بجميع هذه الوجوه والمنافع.... فالخط المجرد
كماله أن تكون دلالته واضحة بإبانة حروفه المتواضعة
وإجادة وضعها ورسمها كل واحد على حدة متميز
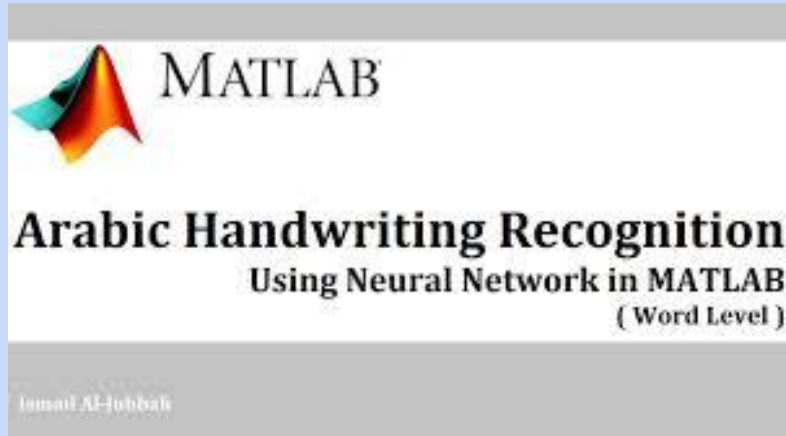عن الاخر إلّا ما اصطلح عليه الكتاب من إيصال حرف

Writing is the outlining and shaping of letters to
indicate audible words which, in turn, indicate
what is in the soul. It comes second after oral
expression. It is a noble craft, since it is one of the
special qualities of man by which he distinguishes
himself from the animals. Furthermore, it reveals
what is in (people's) minds. It enables the intention
(of a person) to be carried to distant places…. It
enables (people) to become acquainted with science,
learning, with the books of the ancients, and with
the sciences and information written down by them.
Because of all these useful aspects, writing is a

# Architecture Design

# Demonstrate Image Parser

# Architecture Design

# *Limitations and Assumption*

**Limitations:**
- Small time and memory constraints (mobile computer)
- Limited connectivity to central database
- Limited central intelligence about documents in other language

**Assume:**
- Documents can be parsed into words
- Words can be mapped to their English word/words (Dictionary lookup)

وأشكال حرفية ⟶ Writing is

# *Classifier Solution*

**Key Idea:**
- Convert document to a simple representation in native language
  - Represent documents as bag of words (BOW) histograms

- This allows train a classifier(s) using corpus of labeled documents in native language **(HUGE CORPUS FROM NATIVE INTELLIGENCE)**
- Classifiers quickly tell what kind of document it is like.
  - Optionally, also classify as "Important for further detailed translation"

**Benefits:**
- Simple classifiers on BOW histograms are just linear feature vectors
  - Low memory, ultra fast decisions
- As more intel is gained, the decision rule(s) can be updated

# User Interface

Command Center with Interpreter

Mobile App

Image

Image Parser

Translate

Search DB

Cache

Elevate?

Yes

No

Ignore

# *Summary*

- The challenge with translating Arabic has to do with image parsing.

- Presented an end-to-end solution with options for Image Parsing, Data Classification and decision making.
  - Google Search concept to find related documents
  - Neural Network techniques
  - Machine Learning techniques

- Questions?

# List of Resources

- [http://arxiv.org/pdf/1206.1518.pdf](http://arxiv.org/pdf/1206.1518.pdf) - research paper for recognizing Arabic characters
- [http://www.amara.org/en/videos/cI4VRKdprrGv/info/arabic-handwriting-recognition-using-neural-network-in-matlab-word-level-user-manual/](http://www.amara.org/en/videos/cI4VRKdprrGv/info/arabic-handwriting-recognition-using-neural-network-in-matlab-word-level-user-manual/) - video tutorial to use Matlab's Arabic Handwriting Recognition feature
- [https://console.ng.bluemix.net/catalog/services/visual-recognition/](https://console.ng.bluemix.net/catalog/services/visual-recognition/) - specs for Watson's Visual Recognition service
-

# Backup Slides

# Language Recognition

- There are thousands of languages spoken on Earth, but the top ten languages cover most speech.

- In field operations, English would be expected to be the default language in only a small percentage of situations.

- Arabic is spoken by

# Modern Standard Arabic

- Challenges of translating content, and also obtaining the right context.

# Simple Yet Complicated

# *What Classifier?*

- Classifier 1: Is it important for detailed translation by human?
- Classifier 2: What type of document is it like?/Which documents that I know are like it?

- You want:
  - Decision confidence
  - Fast lightweight decision
  - Small data decision update
    - **Nearest neighbor:** Distance implicitly tells you confidence about what type of document it is
    - **Relevance vector machine (RVM):** Gives you the category as well as explicit confidence score [0,1]
      Simple weight vector **(super small memory)**

# Limitations

**Hidden assumption:**

- BOW for word -> word translation will give proper representation for classifying using our English corpus.

**If not:**

- Learn an intermediate mapping that re-weights other language BOW to native language BOW.
  - Simple least squares reweighting of histogram
  - Cross correlational analysis