
MADCAT Industry Day



- Joseph Olive
joseph.olive@dar
pa.mil

Announcement



DARPA seeks strong, responsive proposals from well-qualified sources for a new language technology program called Multilingual Automatic Document Classification Analysis and Translation (MADCAT). The goal of this program is to automatically convert foreign language text images into English transcripts, thus eliminating the need for linguists and analysts while automatically providing relevant, distilled actionable information to military command and personnel in a timely fashion.



Goal

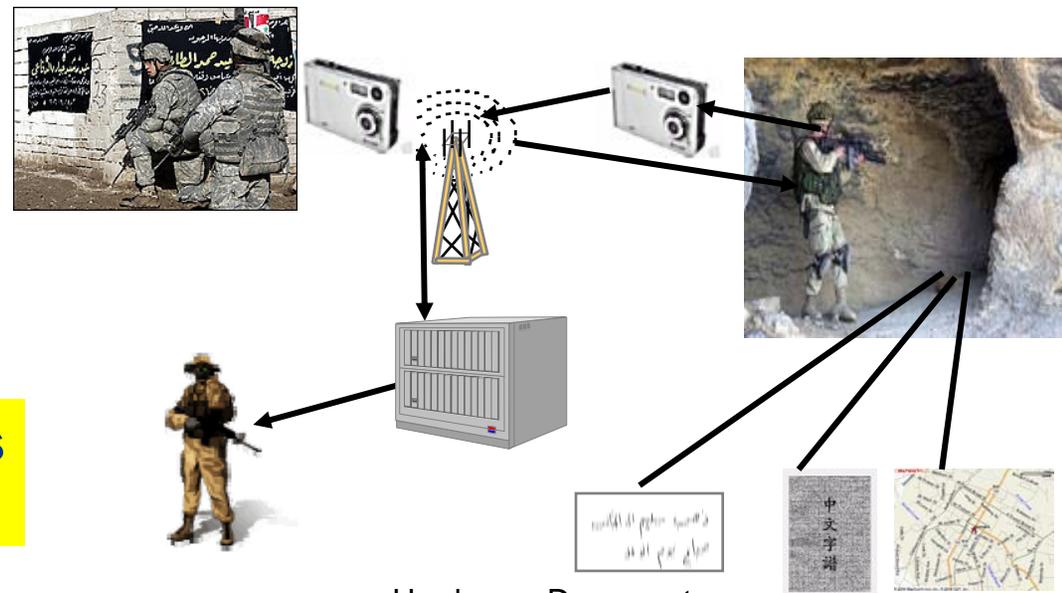


Provide soldiers in the field relevant, distilled, actionable information automatically in a timely fashion from all printed or hand-written Sources in multiple languages

- Strategic Plans, maps
- Financial Records
- Graffiti, propaganda
- Manuals



**Situational awareness
In urban warfare**



Hard-copy Documents



Why is this technology important?



Is this graffiti or is it important?



Necessary for situation awareness – especially for urban warfare

Example Input - Printed



MBC قناة

٦,٠٠ نشرة الأخبار	٧,١٥ ديجيمون
٦,٣٠ من سيربح المليون	٨,١٥ جوني كويست
٨,٠٠ لا تفهمونا غلط	٨,٥٠ أجواء وأطباق
٩,٣٠ مشاهد مثيرة	٩,١٥ صلاة الجمعة - مباشر
١٠,٠١ الأسبوع السياسي	من مكة المكرمة
١١,٠٠ مسلسل عربي	١٠,٠٥ الإفتاء - مباشر
١٢,٣٠ صلالة ٢٠٠١	١١,٣٠ فيلم عربي
١٢,٣٥ بيسي كورة	١,٠٠ الأخبار
١,٠٠ صلالة ٢٠٠١	١,٣٠ برنامج رياضي
٢,٠٠ فيلم أجنبي	١,٥٥ الحياة البرية
٣,٠٠ ديجيمون	٣,٠٠ موجز الأنباء
٣,٢٥ تيمون وبومبا	٣,٠٥ دوري مستر كارد
٣,٥٠ مهووس فيكي	٣,٣٥ حوار الأسبوع
٥,٠٠ عيون الحب	٤,٠٠ موجز الأنباء
٥,٤٥ منوعات غنائية	٤,٠٥ مفكرة المراسل
٦,٠٠ مسلسل عربي	٤,٣٠ الأسبوع السياسي

يستولى على النصيب الأكبر من الصيد، ويشير عليهم ابن السيد البلطي أقدم الصيادين بالعمل معا كشركة للصيد إلا أنهم يتمكنوا من شراء مركب ينافسون به عبد الموجود.

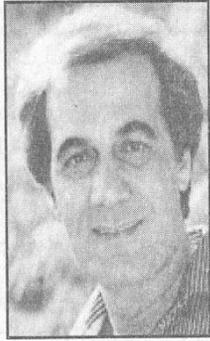
وسهير المرشدي ومديحة حمدي وعبد الرحمن أبو زهرة وإخراج توفيق صالح ويدير حول جماعة من الصيادين يفاجأون بزميلهم عبد الموجود وقد اشترى مركبا كبيرا يعمل بالبخار وبذلك

..و«الثلاثة يحبونها» و«الثار» فيلمان عربيان على الفضائية الثانية

● تعرض القناة الفضائية الثانية الساعة ٢,٠٠ ظهرا بتوقيت القاهرة فيلم (الثلاثة يحبونها) بطولة سعاد نصر وحسن يوسف ويوسف شعبان ويوسف فخر الدين وناهد شريف وإخراج محمود نو الفقار وتدور أحداثه حول الفتاة (إيمان) الموظفة بإحدى الشركات والتي تسيء فهم الحرية وتعيش ثلاث قصص في وقت واحد مع (عصام) مديرها والذي يطلبها للزواج لكنها تصمم على مواصلة دراستها الجامعية، و(كمال) الشاب خفيف الظل و(عادل) الذي ينتقد تصرفاتها المستهترية ● وتعرض الساعة ١,٠٠ بعد منتصف الليل بتوقيت القاهرة فيلم (الثار) بطولة يسرا ومحمود يس وحمدى الوزير وفاروق يوسف وإخراج محمد خان وفيه يقوم أربعة أشخاص باختطاف ندا من زوجها أحمد الذي يلتقط رقم السيارة ويبلغ الشرطة، ويصمم على الانتقام منهم بنفسه ويقوم بقتل ثلاثة أشخاص ويتم القبض عليه قبل ارتكابه الجريمة الرابعة.

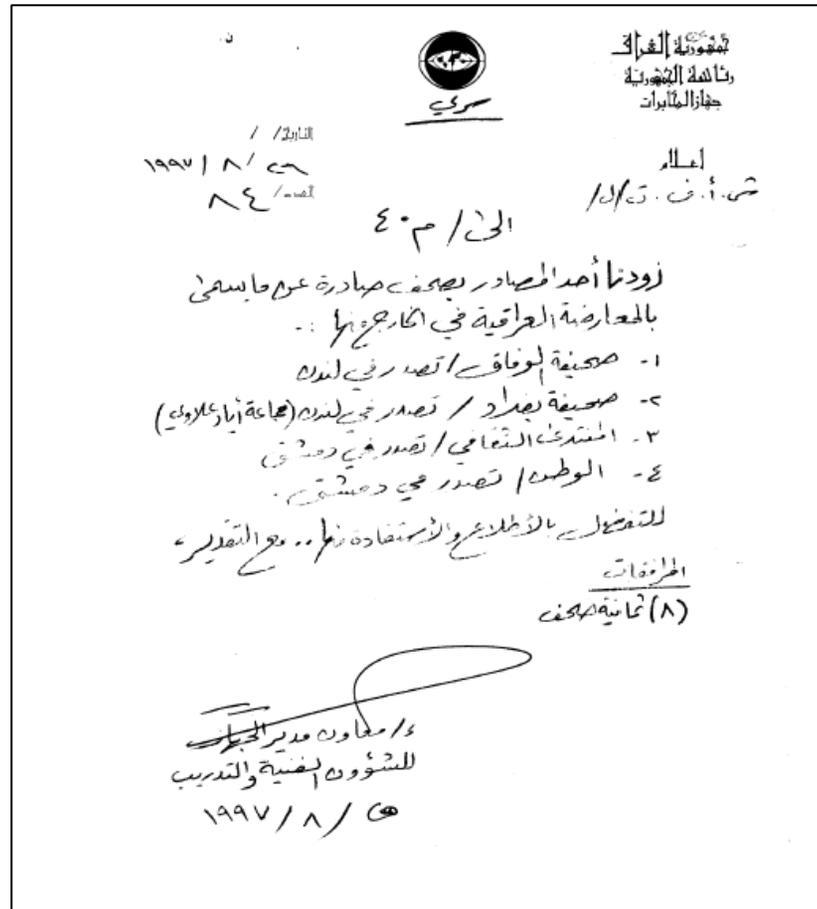
قناة المستقبل

١٣,٣٠ نادى ديزنى	٤,٠٠ كارتون
١٥,٠٠ كميس	٧,٠٠ عالم الصباح
١٦,٠٠ شبكة	٩,٠٠ صلاة الجمعة
١٧,٠٠ الأخبار	١٠,٠٠ المطبخ
١٧,٣٠ ميشو شو	١١,٠٠ موجز الأخبار
١٨,٣٠ فيلم أجنبي	١١,٠٥ فقرة إخبارية
٢٠,٣٠ نص دقيقة	١١,٣٠ فانيسست هوم فيديو
٢١,٣٠ فيلم عربي	١٢,٠٠ الأخبار
٢٣,٣٠ مسلسل عربي	١٣,٠٠ نشرة الأخبار



محمود يس

Example Input - Handwritten



Example Input - Mixed



بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

جمهورية العراق
الجمهورية العراقية
الجمهورية العراقية

امانة سر الحرس الجمهوري

مصري وشقيقي

رئاسة الجمهورية / المكتب الخاص
السكرتير :

م/ الاسرى الكويتيين

تنفيذاً لأمر السيد الرئيس القائد مدام حسين {مخطة اللذ و رعاه} .
مستند ما جاء بقرار مجلس قيادة الثورة في يوم الجمعة
بتاريخ ٣٠٠٣ / ٣ / ٤ .

بنقل كافة الاسرى الكويتيين / البالغين عددهم {٤٤٨} - اسيراً كويتيين .
والمودعين / في سجن {الغمام الاغر} .
المخابرات / المركز العام وسجن كاظمة ب {الكاضمة} .
لجمعهم دروع بشرية في كافة المواقع التي يتوقع ضربها من قبل الامريكان
المعتدين / من مواقع الاتصالات والوزارات الاساسية / الاذاعة والتلفزيون /
مستشفيات التصنيع العسكري /
وكافة المواقع يتوقع الاعتداء عليها من قبل العدوانيين المجرمين الانجلي
امريكان .
تناط مهمة نقلهم بالتنسيق بين كل من :
مديرية جهاز المخابرات
رئاسة اركان الحرس الجمهوري
وباشراف مباشر من قبل جهاز الامن الخاص / امن الجهاز

٢٠٠٣ / ٣ / ٤

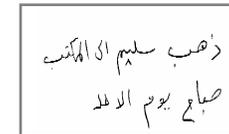
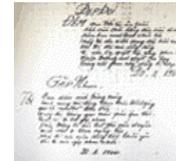
تصميم مدام حسين
المشرف
على امانة سر الحرس الجمهوري

مصورة منه الى :
مديرية جهاز المخابرات / مكتب مدير الجهاز
رئاسة اركان الحرس الجمهوري / مكتب رئيس الأركان

Technology Needs

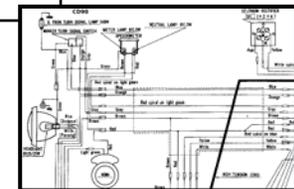


- Create a framework for integration
- Language/Script ID
- Document Type Classification, Segmentation and Labeling
- Zone Analysis and Interpretation
- Robust Multilingual OCR
- Machine Translation
- Distillation



Dr. Joseph Olive
 DARPA
 3701 N. Fairfax, Rm 758
 Arlington, VA 22203
 April, 6 2006

Depart	Arrive	Airline
EWR 10:00AM	SFO 12:30PM	AA
SFO 4:10PM	LAX 5:05PM	UA



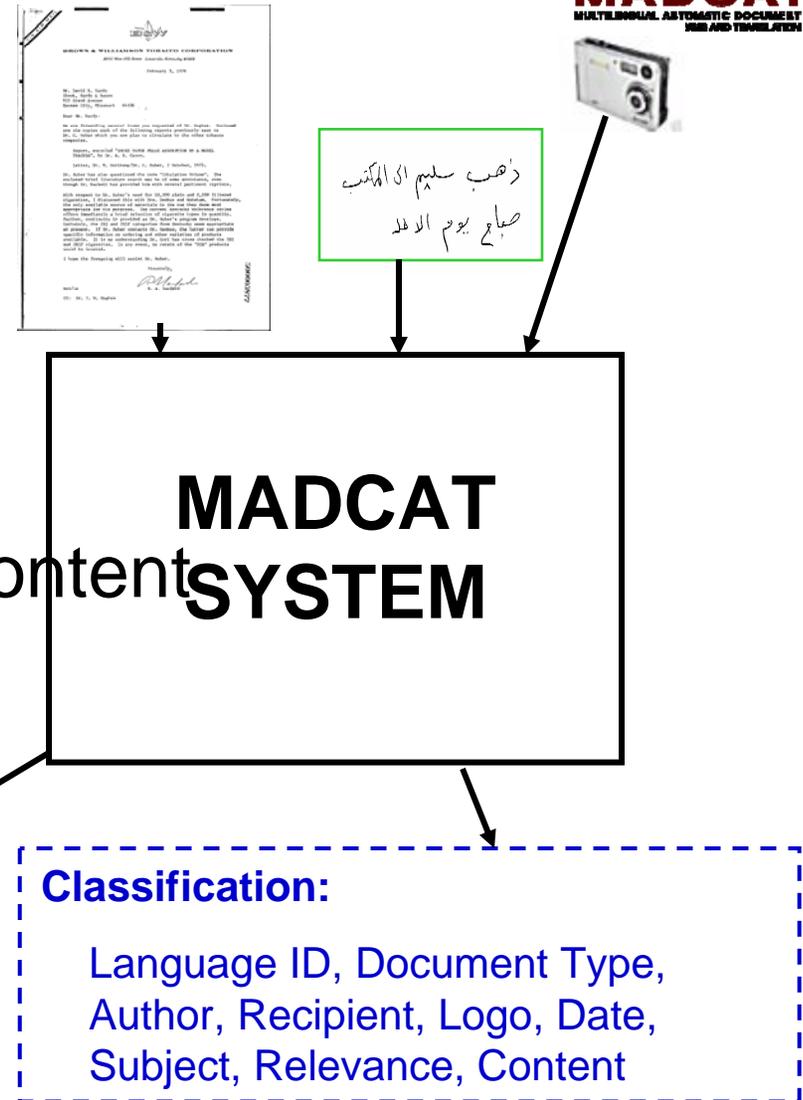
GALE



Deliverables



- Translation
- Classification
 - Author
 - Recipient
 - Names, entities, events content
 - Topics
 - Logos
 - Language ID



Example Output - Printed



English →

Table →

Logo →

Reverse Contrast →

قناة MBC

٦,٠٠ نشرة الأخبار
٦,٣٠ من سيربح المليون
٨,٠٠ لا تفهمونا غلط
٩,٣٠ مشاهد مثيرة
١٠,٠٠ الأسبوع السياسي
١١,٠٠ مسلسل عربي
١٢,٣٠ صلاة ٢٠٠١
١٢,٣٥ بيبيس كورة
١,٠٠ صلاة ٢٠٠١
٢,٠٠ فيلم اجنبي
٣,٠٠ ديجيمون
٣,٢٥ نيمون ويوميا
٣,٥٠ مهبوس فيكي
٥,٠٠ عيون الحب
٥,٤٥ منوعات غنائية
٦,٠٠ مسلسل عربي

٧,١٥ ديجيمون
٨,١٥ جوني كويست
٨,٥٠ اجواء واطباق
٩,١٥ صلاة الجمعة - مباشر
من مكة المكرمة
١٠,٠٥ الإقفاء - مباشر
١١,٣٠ فيلم عربي
١١,٠٠ الأخبار
١,٣٠ برنامج رياضي
١,٥٥ الحياة البرية
٣,٠٠ موجز الأنباء
٣,٠٥ دوري مستر كارد
٣,٣٥ حوار الأسبوع
٤,٠٠ موجز الأنباء
٤,٠٥ فكرة المراسل
٤,٣٠ الأسبوع السياسي

قناة المستقبل

٤,٠٠ كارتون
٧,٠٠ عالم الصباح
٩,٠٠ صلاة الجمعة
١٠,٠٠ المطبخ
١١,٠٠ موجز الأخبار
١١,٠٥ فقرة إخبارية
١١,٣٠ فانيست هوم فيديو
١٢,٠٠ الأخبار
١٣,٠٠ نشرة الأخبار

يستولى على النصيب الأكبر من الصيد، ويشير عليهم ابن السيد البلطي أقدم الصيادين بالعمل معا كشركة للصيد إلا أنهم يتمكنوا من شراء مركب ينافسون به عبد الموجود.

وسهبر المرشدي ومدبحة حمدي وعبد الرحمن أبو زهرة وإخراج توفيق صالح ويدور حول جماعة من الصيادين يفتاجون بزميلهم عبد الموجود وقد اشترى مركبا كبيرا يعمل بالبخار وبذلك

.. والثلاثة يحبونها» و«الثر»
فيلمان عربيان على الفضائية الثانية

● تعرض القناة الفضائية الثانية الساعة ٢,٠٠ ظهرا بتوقيت القاهرة فيلم (الثلاثة يحبونها) بطولة سعاد نصر وحسن يوسف ويوسف شعبان ويوسف فخر الدين وناهد شريف وإخراج محمود نو الفكار وتدور أحداثه حول الفتاة (امان) الموظفة بإحدى الشركات والتي تسيء فهم الحرية وتعيش ثلاث قصص في وقت واحد مع (عصام) مديرها والذي يطلبها للزواج لكنها تصمم على مواصلة دراستها الجامعية، و(كمال) الشاب خفيف الظل و(عادل) الذي يبتعد تصرفاتها المستهتره ● وتعرض الساعة ١,٠٠ بعد منتصف الليل بتوقيت القاهرة فيلم (الثر) بطولة يسرا ومحمود يس وحمدي الوزير وفاروق يوسف وإخراج محمد خان وفيه يقوم أربعة أشخاص باختطاف ندا من زوجها احمد الذي يلتقط رقم السيارة ويبلغ الشرطة، ويصمم على الانتقام منهم بنفسه ويقوم بقتل ثلاثة أشخاص ويتم القبض عليه قبل ارتكابه الجريمة الرابعة.

Two-column text (Arabic)

Text Headline (Arabic)

Graphic (Face)

Caption (Name)

Single column text

Channel MBC

57:1 Djimon	6:00 Who wants to be a millionaire
8:15 Johnny Quest	8:00 Don't misunderstand us
8:50 Points of Views	
9:15 Friday Prayer – Live	
from Mecca Al-Moukrama	
10:05 Religious Decree - Live	
1:00 News	

And Suheer Al-Murshdi and Madeehah Hamdi and Abed Al-Rahman Abu Zahra and directed by Tawfeeq Saleh it is surprised by their mate Abed Al-Mawjoud when he buys a big boat that works on steam

and in return he gains the bigger share of the fishing. Ibin Al-Sayed Al-Balti, who is the oldest fishermen, suggests to them to work together as a fishing company so that they can buy a boat to compete with Abed Al-Mawjoud.

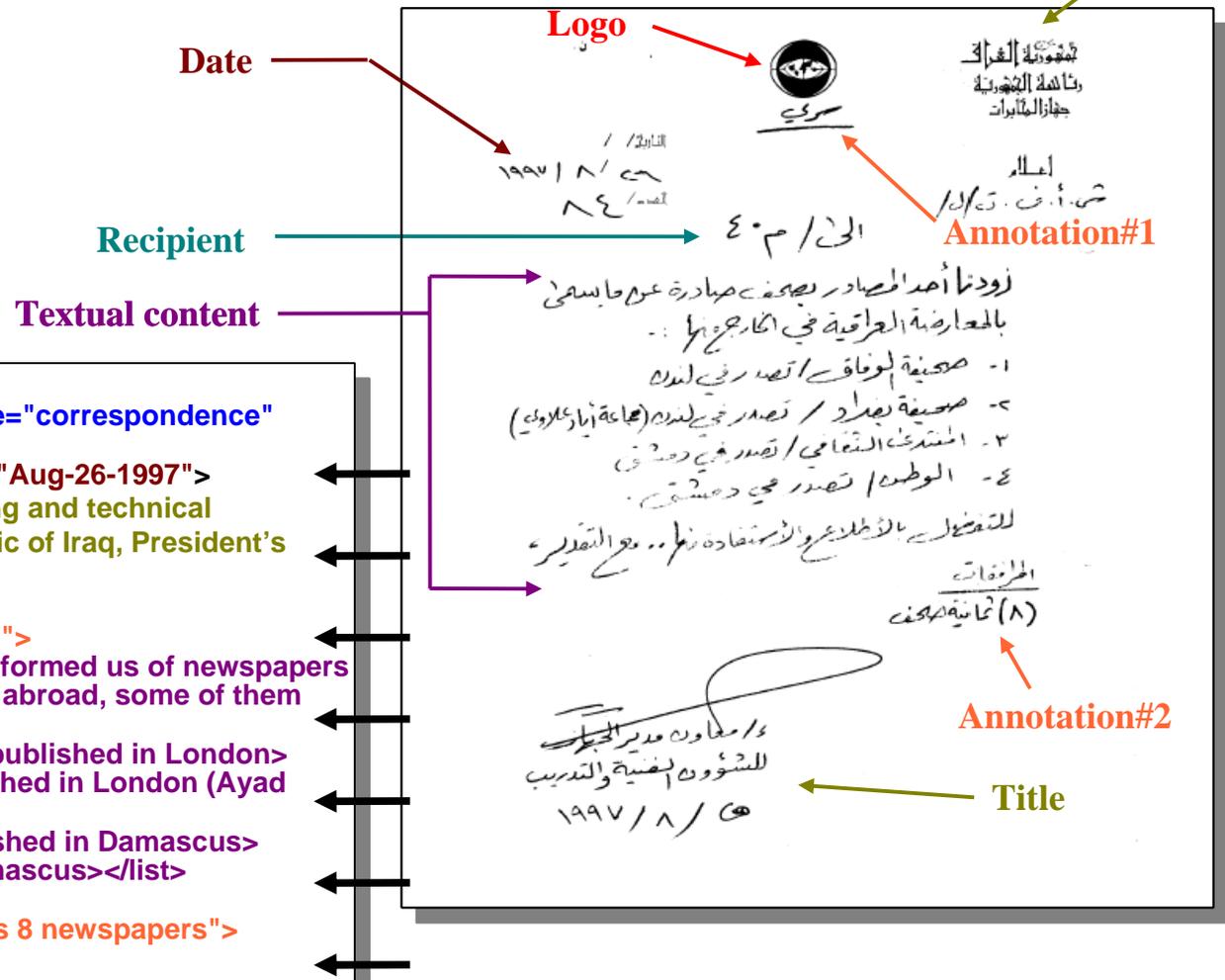
"The Three Love Her" and
"The Retaliation"

two Arabic films on the Al-Fathaeia Al-Thania (Satellite 2)

Example Output - Handwritten



Address



```

<xml>
<meta ID="BIAP-2003-000090.pdf" genre="correspondence"
  language="Arabic" >
<date obtained="Jun-10-2005" created="Aug-26-1997">
<author title="Assistant manager training and technical
  Department" address="The Republic of Iraq, President's
  Office, Intelligence Department">
<recipient name="m40">
<annotation#1 translation="Confidential">
<text translation="One of the sources informed us of newspapers
  that are so called the Iraqi defiance abroad, some of them
  are:"
  <list><item 1-AI Wafaa newspaper published in London>
  <item 2-Baghdad newspaper published in London (Ayad
  Alawi group).>
  <item 3-AI Muftada al Thakafi published in Damascus>
  <item 4-AI Wattan published in Damascus></list>
  "for your information. Regards">
<annotation#2 translation="Attachments 8 newspapers">
</xml>
    
```



Example Output - Mixed



Date →

Recipient →

Textual content →

Logo →

Signature →

Title →

```

<xml>
<meta ID="CMPC-2003-012666.tif" genre="correspondence"
  language="Arabic" >
<date obtained="Mar-22-2006" created="Mar-14-2003">
<author name="Qusai Saddam Hussein"
  title="Supervisor of the Republican Guard Secretariat">
<recipient name="Presidential Office/ Special Office">
...
<text translation="Regarding the execution of Mr. President,
  Commander Saddam Hussein's (God protect him) orders,
  according to the decision of the Revolutionary Command
  Council on Friday, March 4, 2003.

  Transfer all Kuwaiti POW's / a total of 448 captured Kuwaitis
  who are located at the Al-Nida Al-Agher Prison and the
  Intelligence / General Center and Kazema Prison in Al-
  Kazema, to make them human shields at all locations that
  are expected to be attacked by the American aggressors.">
...
</xml>
  
```

Program Requirements



- Four Tasks
 - Printed documents
 - Handwritten and mixed documents
 - Data
 - Evaluation
- Arabic Only
 - Surprise language in later phases
- Only English output evaluated

Evaluation paradigms



- Evaluating Full Translation
 - A set of documents in Arabic (20K words per language)
 - Create a “gold standard”
 - Human editing to match meaning of machine output and “gold standard”
 - Accuracy is $1 - (\text{no. of edits}) / (\text{no. of words})$
- Evaluating Metadata (Logos, Authors, Recipient, Date)
 - A set of documents not seen by the contractors (10K pages)
 - Manually establish ground truth
 - Machine encoded metadata
 - Establish ROC curves and EER points

Metrics



Phase	Translation from Machine Print		Translation from Handwriting		Classification
	% Accuracy	% of documents	% Accuracy	% of documents	% Accuracy
Baseline	65%	70%	2%	50%	60%
Phase 1	75%	80%	40%	70%	75%
Phase 2	85%	85%	60%	80%	85%
Phase 3	90%	85%	75%	85%	95%
Phase 4	95%	90%	85%	90%	Completed
Phase 5	95%	95%	90%	90%	Completed