

# Welche Daten werden benötigt um CEF.AT zu optimieren?

**Bartholomäus Wloka (Universität Wien)**  
Basierend auf der Präsentation von Dr. Khalid Choukri  
(Evaluations and Language resource Distribution Agency)

- Im Vordergrund: ***data-driven paradigm***
  - MÜ Systeme brauchen Daten
  - Fokus für ELRC: Daten in allen Sprachen (EU/CEF)
- Daten, d.h. Sprachressourcen werden produziert aus:
  - Dokumenten & jeglichen Textsammlungen
  - Eine Hilfe von Ihrer Seite ist wichtig

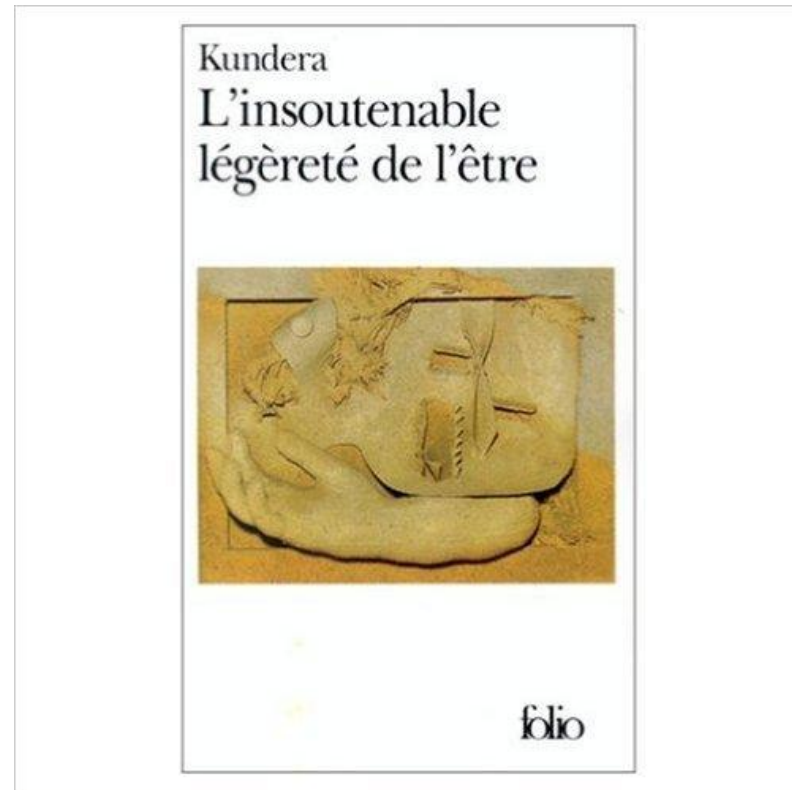
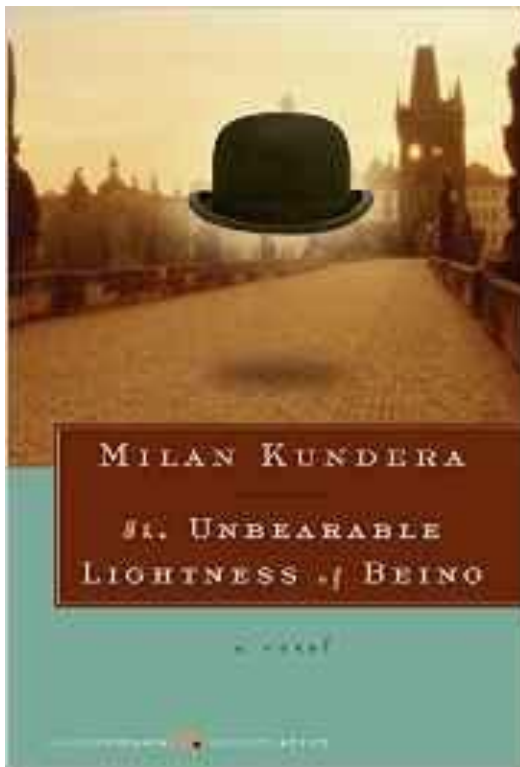


- Alles was Wörter und Sätze beinhaltet, v.A. in mehreren Sprachen
  - Berichte
  - Vorträge
  - Webinhalt
  - Broschüren
  - ...
- Wird zu: *Bag of Words*

# Was sind Daten für MÜ?



wiseGEEK

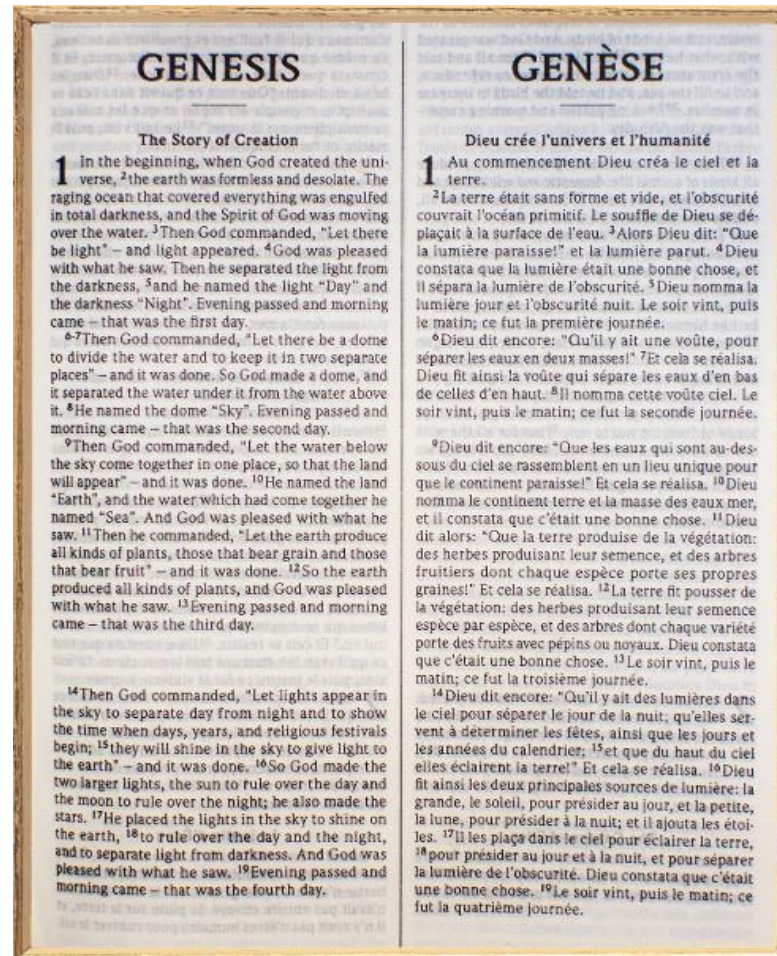




English



French



## English

Telecommunication occurs when the exchange of information between two or more entities (communication) includes the use of technology.

Communication technology uses channels to transmit information (as electrical signals), either over a physical medium (such as signal cables), or in the form of electromagnetic waves.

The word is often used in its plural form, telecommunications, because it involves many different technologies.

## Greek

Με τον γενικό όρο τηλεπικοινωνίες, (telecommunications), χαρακτηρίζεται η κάθε μορφής ενσύρματη ή ασύρματη, ηλεκτρομαγνητική, ηλεκτρική, κ.λπ., ακουστική και οπτική επικοινωνία που πραγματοποιείται ανεξαρτήτως απόστασης.

Στους σύγχρονους καιρούς, αυτή η διαδικασία σχεδόν πάντα περιλαμβάνει την αποστολή ηλεκτρομαγνητικών κυμάτων ή ηλεκτρικών σημάτων από κατάλληλες ηλεκτρονικές συσκευές, όπως το τηλέφωνο ή ο ασύρματος, αλλά παλαιότερα περιελάμβανε τη χρήση ακουστικών σημάτων, όπως τυμπάνων, ή οπτικών, όπως ο σηματοφόρος καπνός ή η λάμψη της φωτιάς.

## Spanish

Una telecomunicación es toda transmisión y recepción de señales de cualquier naturaleza, típicamente electromagnéticas, que contengan signos, sonidos, imágenes o, en definitiva, cualquier tipo de información que se desee comunicar a cierta distancia.

Por metonimia, también se denomina telecomunicación (o telecomunicaciones, indistintamente) a la disciplina que estudia, diseña, desarrolla y explota aquellos sistemas que permiten dichas comunicaciones; de forma análoga, la ingeniería de telecomunicaciones resuelve los problemas técnicos asociados a esta disciplina.

**Source:** First sentences of articles for Telecommunications in the English, Greek and Spanish Wikipedias

**German page is slightly different but these are (never) translations of one source!!**



The Vikings were Scandinavian seafarers who lived in the ninth, tenth, and the beginning of the eleventh century, which is known as the Viking era. The Vikings were heathens and did not become Christian until around the year 1000. Their own gods were called the Æsir, and offerings were made to them at the blot, a kind of religious sacrificial holiday.

Four of these gods were Tyr (or Tiwaz), Odin (or Wotan), Thor, and Frigga, who have given their names to four of the days of the week: Tuesday, Wednesday, Thursday and Friday. The months had their own names as well, but now the Scandinavians use the Roman names for the months: January, February, March etc.

Many Vikings sailed out into the world in their long-ships, or drekkrar, as far as America and Constantinople. Their ships had relatively flat bottoms, so that they could sail near the coast and up shallow rivers. In the West they met Indians, and in the East they met Arabs. Out in the Atlantic they navigated by the stars, and in the year 1000 Leif Eriksson set foot on American soil, and forty years later, Ingvar the Wide-Traveled reached the southern shore of the Caspian sea. In this way, local kings had contact with lands which lay far away. In large areas of England Danish law held sway; that area was therefore called the Danelaw. In Constantinople, the emperor had a feared bodyguard composed of Vikings. Because of their distinctive axes, they were called "the Axe-bearing Barbarians."

At home the Vikings lived relatively simply. They sowed rye in the fields and kept cows, which gave milk, pigs, for pork, and sheep, for wool. Those who lived along the coasts caught fish. They often lived in long-houses, which could house several families. Three or four brothers, for example, could live with their families together in one big house.

Die Wikinger waren skandinavische Seefahrer, die im 9., 10. und Anfang des 11. Jahrhunderts lebten, auch bekannt als Wikinger-Epoche. Die Wikinger waren Heiden und wurden erst um das Jahr 1000 zu Christen. Ihre eigenen Götter nannten sie Æsir, denen sie am Blot, einem religiösen Opfertag, Gaben darbrachten. Vier dieser Götter waren Tyr (oder Tiwaz), Odin (oder Wotan), Thor und Frigga, nach denen drei Wochentage benannt sind: Dienstag, Donnerstag und Freitag. Auch die Monate hatten ihre eigenen Namen, aber heutzutage benutzen die Skandinavier die römischen Namen für die Monate: Januar, Februar, März etc.

Viele Wikinger segelten in ihren Langschiffen oder Drekkar hinaus in die Welt, bis nach Amerika und Konstantinopel. Ihre Schiffe hatten relativ flache Böden, so daß sie sich damit auch nahe der Küste und in seichten Flüssen bewegen konnten.

Im Westen begegneten sie Indianern und im Osten Arabern. Auf dem Atlantik navigierten sie mit Hilfe der Sterne und im Jahr 1000 setzte Leif Eriksson seinen Fuß auf amerikanischen Boden, und vierzig Jahre später erreichte Ingvar, 'der Weitgereiste', die Südküste des Kaspischen Meeres. Auf diese Weise kamen einheimische Könige in Kontakt mit Ländern, die weit entfernt waren.

In weiten Teilen Englands herrschte dänisches Gesetz. Diese Gebiete wurden deshalb Danelaw genannt. In Konstantinopel hielt sich der Herrscher eine gefürchtete Wikingergarde. Wegen ihrer typischen Streitäxte wurden sie die Axt-tragenden Barbaren genannt.

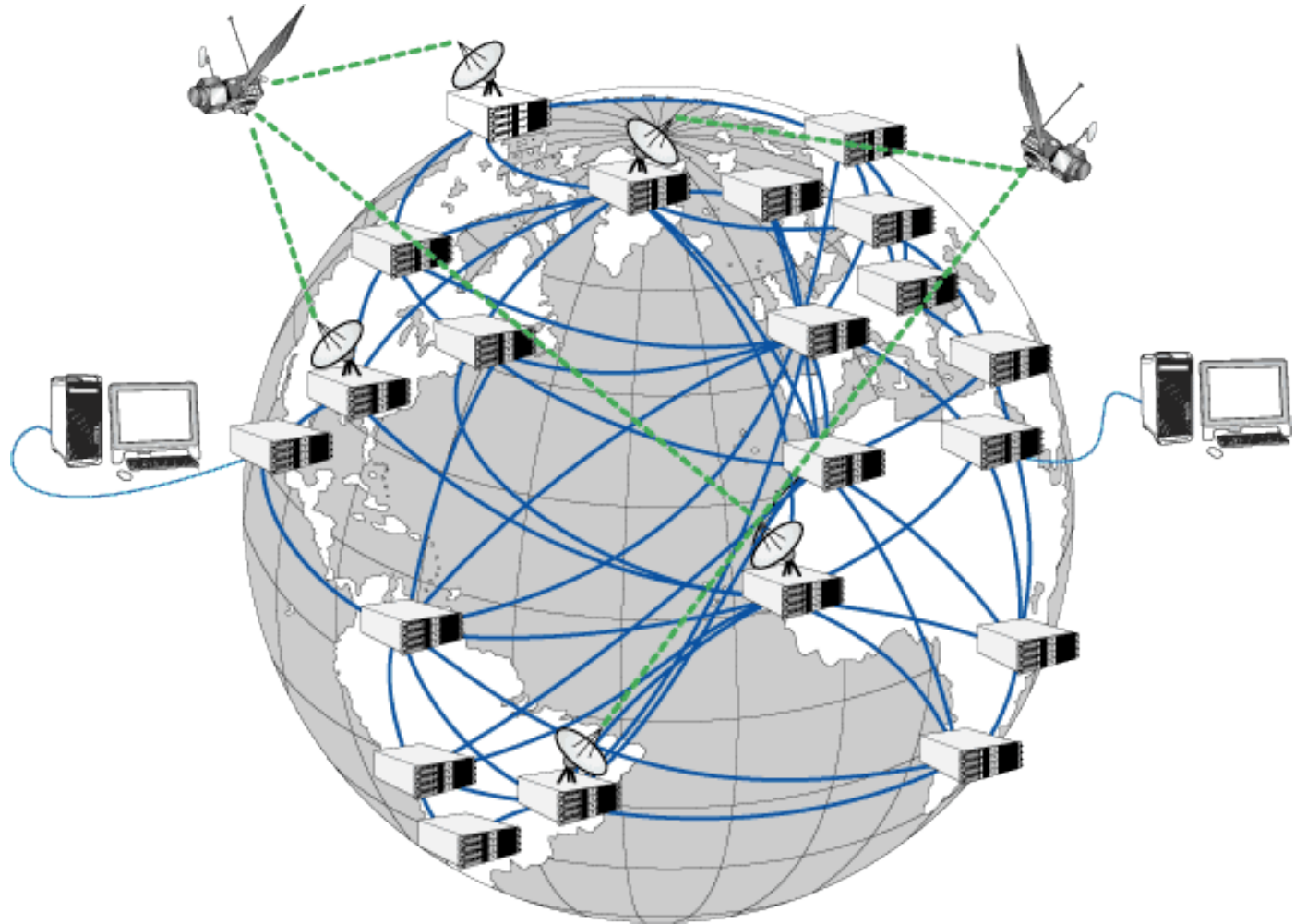
Zu Hause lebten die Wikinger recht einfach. Auf den Feldern kultivierten sie Roggen und sie hielten Kühe, die sie mit Milch versorgten. Schweine hielten sie wegen des Fleisches und Schafe für Wolle. Jene, die an der Küste lebten, fingen Fisch. Die Wikinger wohnten gewöhnlich in Langhäusern, die mehrere Familien beherbergen konnten. Drei oder vier Brüder konnten, zum Beispiel, zusammen mit ihren Familien in einem einzigen großen Haus leben.

highly ...  
previous level in time or space.

| ID   | FR                  | ES                     | EL   |
|------|---------------------|------------------------|--|
| 6905 | abandon scolaire    | abandono escolar       | διακοπή της σχολικής φοίτησης                      |
| 920  | abats               | despojo                | παραπροϊόντα σφαγίων                               |
| 1857 | abattage d'animaux  | sacrificio de animales | σφαγή ζώων   |
| 6621 | abrogation          | derogación             | κατάργηση  |
| 5075 | Abruzzes            | Abruzos                | Αβρουζία<br>συστηματική απουσία από την<br>εργασία |
| 5339 | absentéisme         | absentismo             | εργασία  |
| 5984 | abstentionnisme     | abstencionismo         | αποχή  |
| 2    | abus de confiance   | abuso de confianza     | απιστία  |
| 25   | abus de droit       | abuso de derecho       | κατάχρηση δικαιώματος                              |
|      | abus de pouvoir     | abuso de poder         | κατάχρηση εξουσίας                                 |
|      | accès à l'éducation | acceso a la educación  | πρόσβαση στην εκπαίδευση                           |
|      | accès à l'emploi    | acceso al empleo       | πρόσβαση στην αγορά εργασίας                       |



# Wo sind diese Daten zu finden?



# Welche Formate werden benötigt?





## Dublin Core Metadata Element Set

1. Title
2. Creator
3. Subject
4. Description
5. Publisher
6. Contributor
7. Date
8. Type
9. Format
10. Identifier
11. Source
12. Language
13. Relation
14. Coverage
15. Rights



- Beispiele von Rohdaten (HTML mit Bildern, Tabellen, etc.) die in Sprachrecources umgewandelt werden
  - Auffinden und identifizieren der Daten
  - Erwerb der Daten (z.B. mit einem Crawler)
  - Daten säubern (d.h. entfernen von Standardtexten, Menus, Bildern, HTML Tags, etc. und Formatkonvertierung)
  - Beispiel: *Boilerpipe*
  - Dokumentation der Daten
  - Alignierung der Übersetzungen auf Satzebene.
  - Berechnung der Alignierungsqualität
  - Share

## La France en Allemagne

## Ein Ausgangspunkt sind Webseiten

6 juillet 2015

### Elysée : rencontre sur la Grèce avec Angela Merkel

La chancelière allemande Angela Merkel et le président de la République François Hollande auront, lundi 6 juillet au soir à l'Elysée, un entretien suivi d'un dîner de travail pour évaluer les conséquences du référendum en Grèce.

Cette rencontre s'inscrit dans le cadre de la coopération permanente entre la France et l'Allemagne pour contribuer à une solution durable en (...)



### Réformes en France



### Paris Climat 2015 / COP21 - Pour un accord universel



## Frankreich in Deutschland

6. Juli 2015

### Griechenland: Staatspräsident Hollande empfängt Bundeskanzlerin (...)

Staatspräsident François Hollande empfängt am Abend des 6. Juli Bundeskanzlerin Angela Merkel im Elysée-Palast zu einem Gespräch und einem Arbeitsessen, um die Konsequenzen aus dem Referendum in Griechenland zu erörtern. Das Treffen findet im Rahmen der ständigen Zusammenarbeit zwischen Frankreich und Deutschland mit dem Ziel statt, zu einer dauerhaften Lösung für Griechenland zu (...)



### Reformagenda



### Klimakonferenz Paris 2015







## Griechenland: Staatspräsident Hollande empfängt Bundeskanzlerin Merkel [\[fr\]](#)

Drucken

[Google](#)

[Facebook](#)

[Twitter](#)

Staatspräsident François Hollande empfängt am Abend des 6. Juli Bundeskanzlerin Angela Merkel im Elysée-Palast zu einem Gespräch und einem Arbeitsessen, um die Konsequenzen aus dem Referendum in Griechenland zu erörtern.

Das Treffen findet im Rahmen der ständigen Zusammenarbeit zwischen Frankreich und Deutschland mit dem Ziel statt, zu einer dauerhaften Lösung für Griechenland zu kommen..

Letzte Änderung 06/07/2015

[Seitenanfang](#)

## Elysée : rencontre sur la Grèce avec Angela Merkel [\[de\]](#)

• Imprimer

• [Google](#)

• [Facebook](#)

• [Twitter](#)

• La chancelière allemande Angela Merkel et le président de la République François Hollande auront, lundi 6 juillet au soir à l'Elysée, un entretien suivi d'un diner de travail pour évaluer les conséquences du référendum en Grèce.

• Cette rencontre s'inscrit dans le cadre de la coopération permanente entre la France et l'Allemagne pour contribuer à une solution durable en Grèce.

• Dernière modification : 06/07/2015

• [Haut de page](#)

## ➤ RAW DATA TO PROCESS

## Griechenland: Staatspräsident Hollande **empfängt** Bundeskanzlerin Merkel

François Hollande empfängt am Abend des 6. Juli Bundeskanzlerin Angela Merkel im **Elysée-Palast** zu einem Gespräch und einem Arbeitsessen, um die Konsequenzen aus dem Referendum in Griechenland zu erörtern.

Das Treffen findet im Rahmen der ständigen Zusammenarbeit zwischen Frankreich und Deutschland mit dem Ziel statt, zu einer dauerhaften Lösung für Griechenland zu kommen..

Letzte Änderung 06/07/2015

[Seitenanfang](#)

- **Elysée : **rencontre** sur la Grèce avec Angela Merkel**

- La chancelière allemande Angela Merkel et le président de la République François Hollande auront, lundi 6 juillet au soir à **l'Elysée**, un entretien suivi d'un diner de travail pour évaluer les conséquences du référendum en Grèce.

- Cette rencontre s'inscrit dans le cadre de la coopération permanente entre la France et l'Allemagne pour contribuer à une solution durable en Grèce.

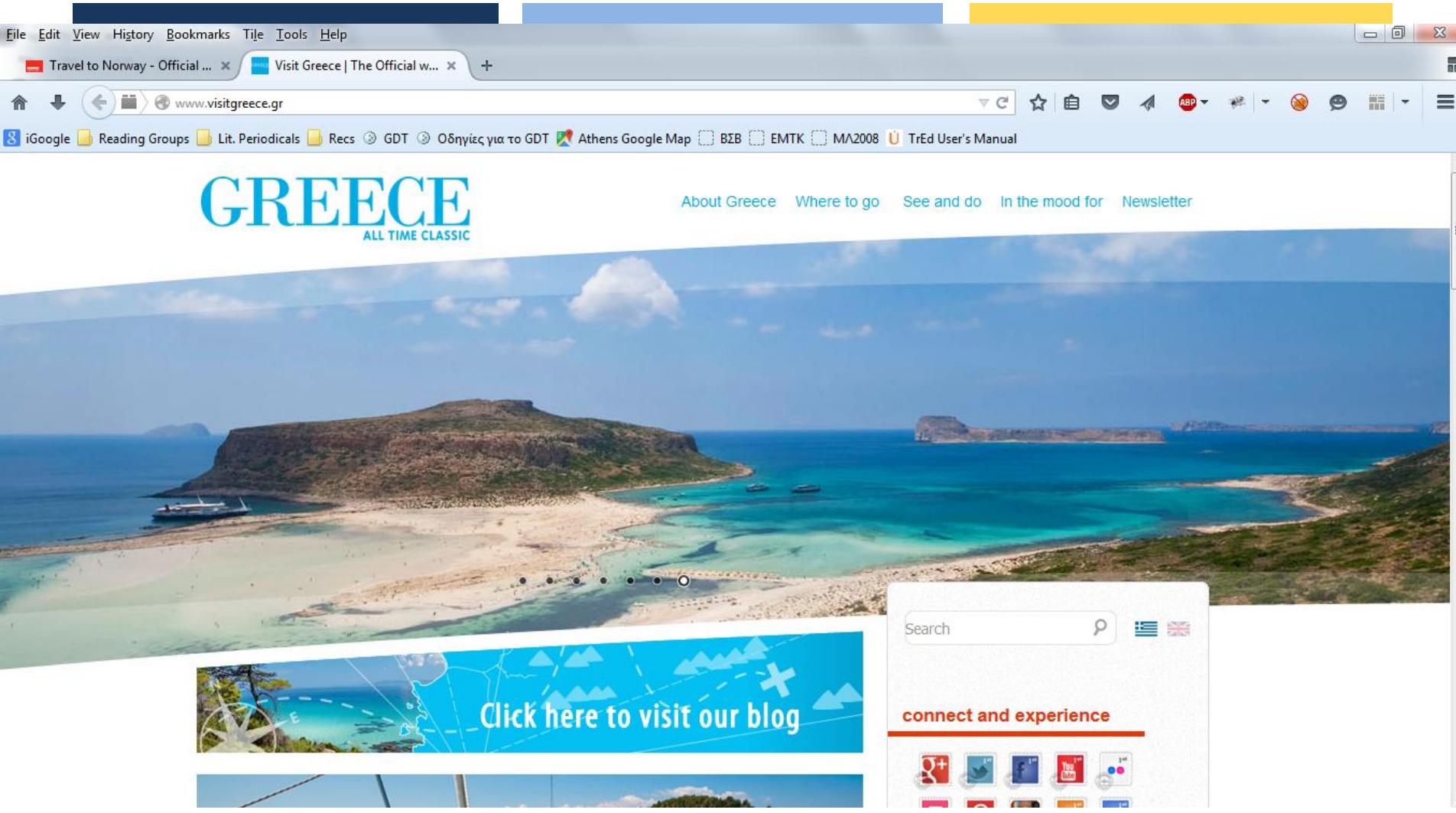
- Dernière modification : 06/07/2015

- [Haut de page](#)

➤ **RAW DATA TO PROCESS**

- Wie kann dieser Prozess automatisiert werden?
- Wir fangen an mit digital verfügbaren Daten
  - OCR kann für Sprachen eingesetzt die weniger in digitaler Form vertreten sind

# Viele Seiten bieten Inhalt in mehreren Sprachen an



File Edit View History Bookmarks Title Tools Help

Travel to Norway - Official ... Visit Greece | The Official w...

www.visitgreece.gr



iGoogle Reading Groups Lit. Periodicals Recs GDT Οδηγίες για το GDT Athens Google Map ΒΣΒ EMTK ΜΑ2008 TrEd User's Manual

# GREECE






ALL TIME CLASSIC

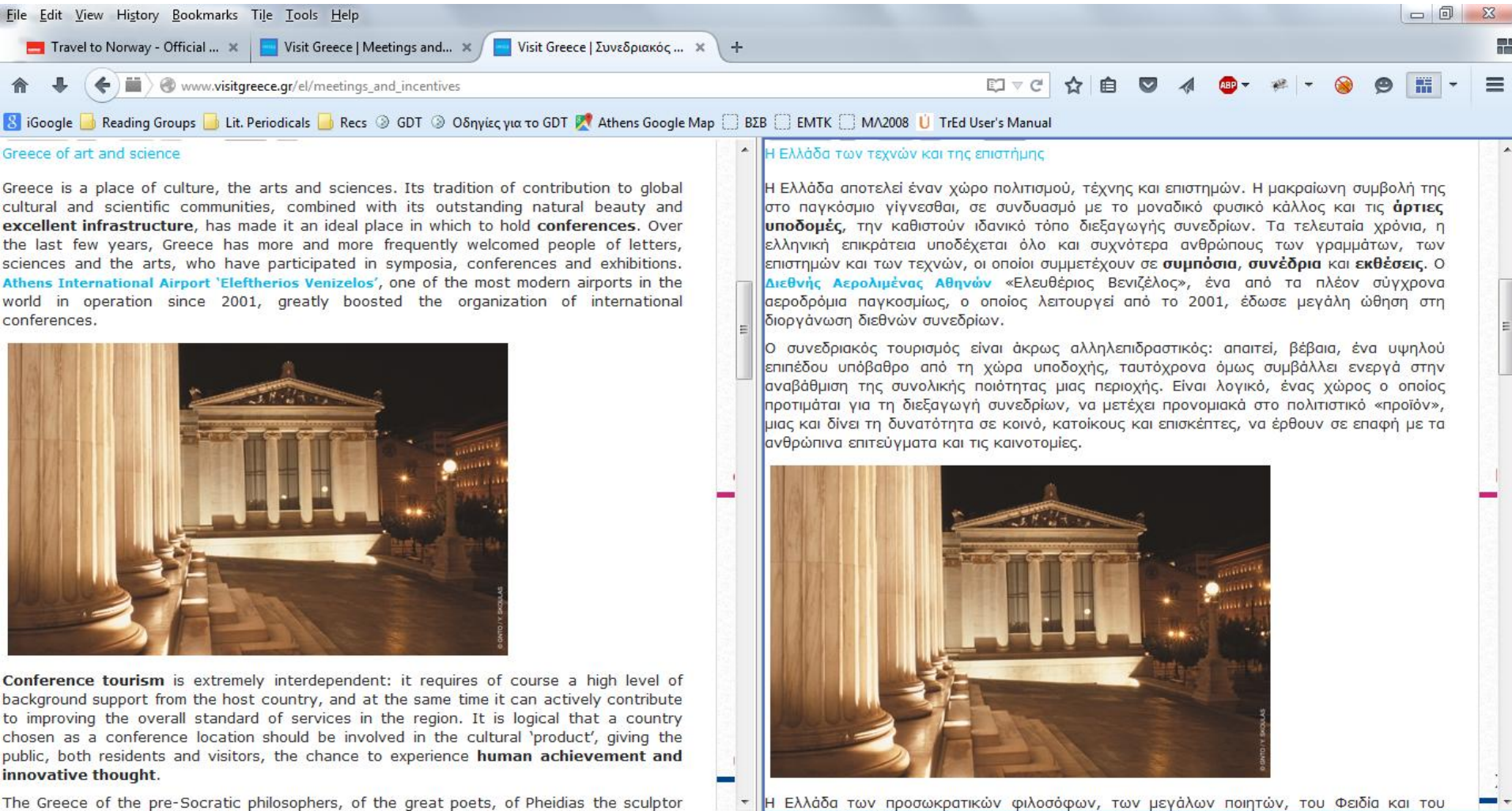
About Greece Where to go See and do In the mood for Newsletter

Click here to visit our blog

Search  

connect and experience



File Edit View History Bookmarks Title Tools Help

Travel to Norway - Official ... x Visit Greece | Meetings and... x Visit Greece | Συνεδριακός ... x +

www.visitgreece.gr/el/meetings\_and\_incentives

iGoogle Reading Groups Lit. Periodicals Recs GDT Οδηγίες για το GDT Athens Google Map ΒΣΒ EMTK ΜΑ2008 TrEd User's Manual

Greece of art and science

Greece is a place of culture, the arts and sciences. Its tradition of contribution to global cultural and scientific communities, combined with its outstanding natural beauty and **excellent infrastructure**, has made it an ideal place in which to hold **conferences**. Over the last few years, Greece has more and more frequently welcomed people of letters, sciences and the arts, who have participated in symposia, conferences and exhibitions. **Athens International Airport 'Eleftherios Venizelos'**, one of the most modern airports in the world in operation since 2001, greatly boosted the organization of international conferences.

**Conference tourism** is extremely interdependent: it requires of course a high level of background support from the host country, and at the same time it can actively contribute to improving the overall standard of services in the region. It is logical that a country chosen as a conference location should be involved in the cultural 'product', giving the public, both residents and visitors, the chance to experience **human achievement and innovative thought**.

The Greece of the pre-Socratic philosophers, of the great poets, of Pheidias the sculptor

Η Ελλάδα των τεχνών και της επιστήμης

Η Ελλάδα αποτελεί έναν χώρο πολιτισμού, τέχνης και επιστημών. Η μακραίωνη συμβολή της στο παγκόσμιο γίνεσθαι, σε συνδυασμό με το μοναδικό φυσικό κάλλος και τις **άρτιες υποδομές**, την καθιστούν ιδανικό τόπο διεξαγωγής συνεδρίων. Τα τελευταία χρόνια, η ελληνική επικράτεια υποδέχεται όλο και συχνότερα ανθρώπους των γραμμάτων, των επιστημών και των τεχνών, οι οποίοι συμμετέχουν σε **συμπόσια, συνέδρια και εκθέσεις**. Ο **Διεθνής Αερολιμένας Αθηνών** «Ελευθέριος Βενιζέλος», ένα από τα πλέον σύγχρονα αεροδρόμια παγκοσμίως, ο οποίος λειτουργεί από το 2001, έδωσε μεγάλη ώθηση στη διοργάνωση διεθνών συνεδρίων.

Ο συνεδριακός τουρισμός είναι άκρως αλληλεπιδραστικός: απαιτεί, βέβαια, ένα υψηλού επιπέδου υπόβαθρο από τη χώρα υποδοχής, ταυτόχρονα όμως συμβάλλει ενεργά στην αναβάθμιση της συνολικής ποιότητας μιας περιοχής. Είναι λογικό, ένας χώρος ο οποίος προτιμάται για τη διεξαγωγή συνεδρίων, να μετέχει προνομιακά στο πολιτιστικό «προϊόν», μιας και δίνει τη δυνατότητα σε κοινό, κατοίκους και επισκέπτες, να έρθουν σε επαφή με τα ανθρώπινα επιτεύγματα και τις καινοτομίες.

Η Ελλάδα των προσωκρατικών φιλοσόφων, των μεγάλων ποιητών, του Φειδία και του

File Edit View History Bookmarks Title Tools Help

Travel to Norway - Official ... x Visit Greece | Meetings and... x Visit Greece | Συνεδριακός ... x Sentence alignment for 103.xml... x +

abumatan.eu/~vpapa/data/EN-EL/crawled\_data/visitgreece\_20150825\_154605/eac25a8b-87cd-4b08-b045-571ccb003af6/xml/1234\_103\_u.tmx.html

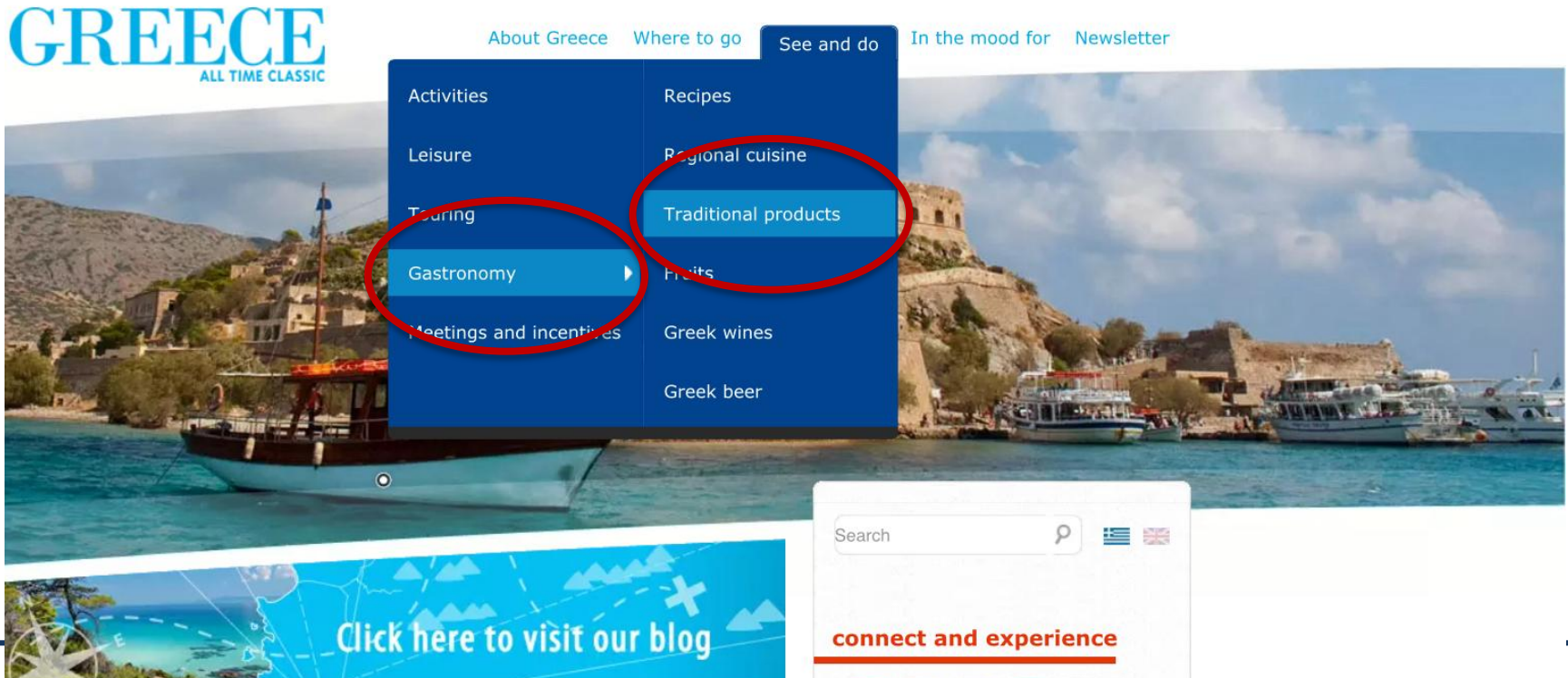
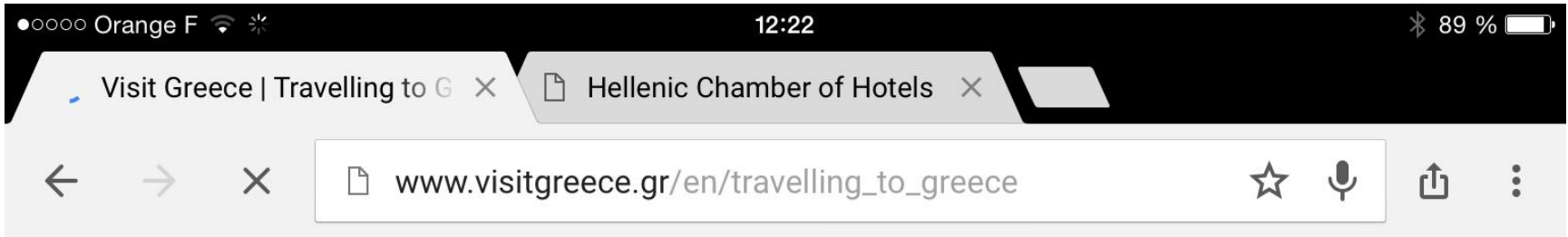
iGoogle Reading Groups Lit. Periodicals Recs GDT Οδηγίες για το GDT Athens Google Map ΒΣΒ EMTK ΜΑ2008 TrEd User's Manual

## Sentence alignment for 103.xml (en) - 1234.xml (el)

| #  | en  | el  |
|----|---|---|
| 1  | Greece of art and science   | Η Ελλάδα των τεχνών και της επιστήμης   |
| 2  | Greece is a place of culture, the arts and sciences.  | Η Ελλάδα αποτελεί έναν χώρο πολιτισμού, τέχνης και επιστημών.   |
| 3  | Its tradition of contribution to global cultural and scientific communities, combined with its outstanding natural beauty and excellent infrastructure, has made it an ideal place in which to hold conferences.  | Η μακράννη συμβολή της στο παγκόσμιο γίνεσθαι, σε συνδυασμό με το μοναδικό φυσικό κάλλος και τις άρτιες υποδομές, την καθιστούν ιδανικό τόπο διεξαγωγής συνεδρίων.  |
| 4  | Over the last few years, Greece has more and more frequently welcomed people of letters, sciences and the arts, who have participated in symposia, conferences and exhibitions.   | Τα τελευταία χρόνια, η ελληνική επικράτεια υποδέχεται όλο και συχνότερα ανθρώπους των γραμμάτων, των επιστημών και των τεχνών, οι οποίοι συμμετέχουν σε συμπόσια, συνέδρια και εκθέσεις.  |
| 5  | Athens International Airport 'Eleftherios Venizelos', one of the most modern airports in the world in operation since 2001, greatly boosted the organization of international conferences.  | Ο Διεθνής Αερολιμένας Αθηνών «Ελευθέριος Βενιζέλος», ένα από τα πλέον σύγχρονα αεροδρόμια παγκοσμίως, ο οποίος λειτουργεί από το 2001, έδωσε μεγάλη ώθηση στη διοργάνωση διεθνών συνεδρίων.   |
| 6  | Conference tourism is extremely interdependent: it requires of course a high level of background support from the host country, and at the same time it can actively contribute to improving the overall standard of services in the region.            | Ο συνεδριακός τουρισμός είναι άκρως αλληλεπιδραστικός: απαιτεί, βέβαια, ένα υψηλού επιπέδου υπόβαθρο από τη χώρα υποδοχής, ταυτόχρονα όμως συμβάλλει ενεργά στην αναβάθμιση της συνολικής ποιότητας μιας περιοχής.  |
| 7  | It is logical that a country chosen as a conference location should be involved in the cultural 'product', giving the public, both residents and visitors, the chance to experience human achievement and innovative thought.                           | Είναι λογικό, ένας χώρος ο οποίος προτιμάται για τη διεξαγωγή συνεδρίων, να μετέχει προνομιακά στο πολιτιστικό «προϊόν», μιας και δίνει τη δυνατότητα σε κοινό, κατοίκους και επισκέπτες, να έρθουν σε επαφή με τα ανθρώπινα επιτεύγματα και τις καινοτομίες. |
| 8  | The Greece of the pre-Socratic philosophers, of the great poets, of Pheidias the sculptor and Asclepius the physician, extends its hospitality and its warmest welcome, honouring people of intellect and creativity, commerce and scientific progress. | Η Ελλάδα των προσωκρατικών φιλοσόφων, των μεγάλων ποιητών, του Φειδία και του Ασκληπιού υποδέχεται φιλόξενα και τιμά τους ανθρώπους του πνεύματος, του εμπορίου και της προόδου.  |
| 9  | Scientific conferences in the land that gave birth to science   | Συνέδρια στη χώρα που γέννησε τις επιστήμες   |
| 10 | Greece has a large number of esteemed scientists, both here in the country and abroad.  | Η Ελλάδα διαθέτει μεγάλο και υψηλής αξίας επιστημονικό δυναμικό, τόσο εντός όσο και εκτός συνόρων.  |
| 11 | Greek scientists, with their inventions, innovations and research work, play a leading part in the international scientific community.  | Οι Έλληνες επιστήμονες, με τις εφευρέσεις, τις καινοτομίες και το ερευνητικό τους έργο πρωταγωνιστούν στη διεθνή επιστημονική κοινότητα.  |
| 12 | Numerous important scientific conferences take place in Greece, reflecting the significance the country places on innovative science.   | Τα επιστημονικά συνέδρια που λαμβάνουν χώρα στην Ελλάδα είναι και πολλά και σημαντικά, αντανακλώντας τη σημασία που δίνει η χώρα στις καινοτόμες επιστήμες.   |
| 13 | Medical, architectural, natural and humanistic scientific conferences enrich Greece's cultural life, and at the same time give participants the opportunity to experience the   | Ιατρικά συνέδρια, αρχιτεκτονικά, φυσικών και ανθρωπιστικών επιστημών, πλουτίζουν την πολιτιστική ζωή της  |

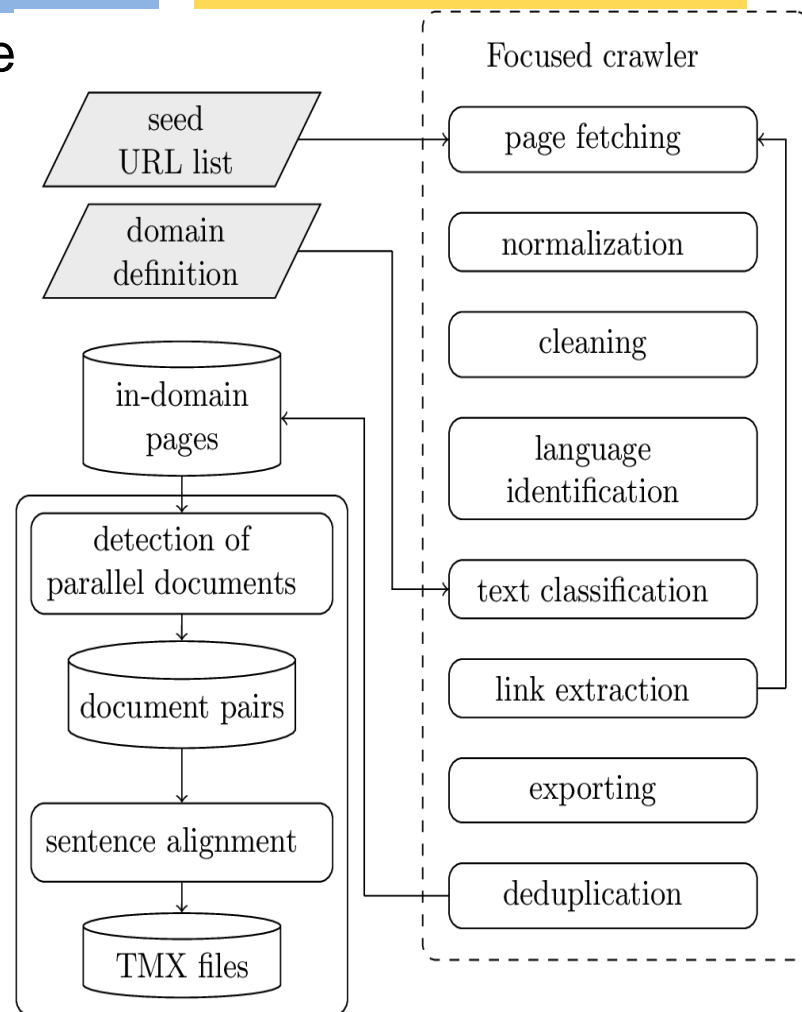
- Wir gehen vor indem wir:
  - Eine Quelle als nützlich identifizieren (gute multilinguale Daten)
  - **Alle weiteren Links abarbeiten**
    - **Analog zum anklicken eines jeden Links (Crawling)**
- Erhalt der Seiten mit der zugehörigen übersetzten Seiten
- Identifikation von Domain und Genre, wenn möglich
- Säuberung der Daten
- Alignierung (auf Dokument-, Abschnitt-, Satzebene)
- Überprüfung der Alignierung
- Use & Share

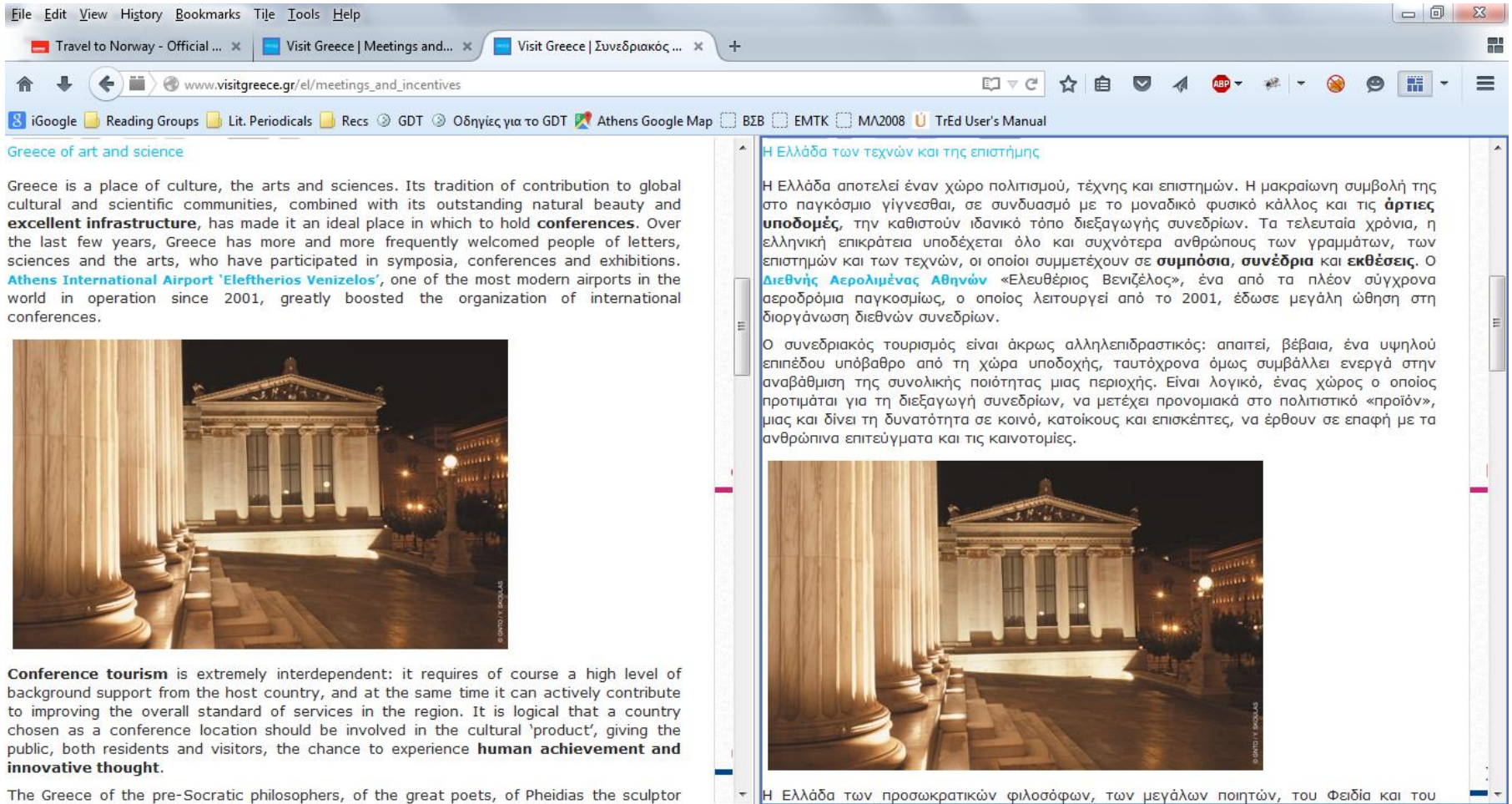
# Automatisches Verfolgen aller referenzierter Links





- Prototyp zum Erwerb von generellen, oder Domain-spezifischen ein- oder mehrsprachigen Korpora
- **Eingabe:**
  - Definition von Domänen (Termini)
  - Seed URL
- Module(open source libraries/toolkits):
  - Page Fetching/Text Extraction
  - Normalization and Metadata Extraction
  - Boilerplate Detection (Boilerpipe)
  - Spracherkennung (> 50 Sprachen)
  - Textklassifikation
  - Entfernen von Duplikaten
  - Erkennen der Dokumentenpaare
  - Satzalignierung (Hunalign and others)
- **Ausgabe**
  - Dokumentenpaare
  - Segmentpaare in TMX Format





File Edit View History Bookmarks Title Tools Help  
 Travel to Norway - Official ... Visit Greece | Meetings and... Visit Greece | Συνεδριακός ...  
 www.visitgreece.gr/el/meetings\_and\_incentives  
 iGoogle Reading Groups Lit. Periodicals Recs GDT Οδηγίες για το GDT Athens Google Map ΒΣΒ EMTK ΜΑ2008 TrEd User's Manual  
 Greece of art and science  
 Greece is a place of culture, the arts and sciences. Its tradition of contribution to global cultural and scientific communities, combined with its outstanding natural beauty and **excellent infrastructure**, has made it an ideal place in which to hold **conferences**. Over the last few years, Greece has more and more frequently welcomed people of letters, sciences and the arts, who have participated in symposia, conferences and exhibitions. **Athens International Airport 'Eleftherios Venizelos'**, one of the most modern airports in the world in operation since 2001, greatly boosted the organization of international conferences.

**Conference tourism** is extremely interdependent: it requires of course a high level of background support from the host country, and at the same time it can actively contribute to improving the overall standard of services in the region. It is logical that a country chosen as a conference location should be involved in the cultural 'product', giving the public, both residents and visitors, the chance to experience **human achievement and innovative thought**.

The Greece of the pre-Socratic philosophers, of the great poets, of Pheidias the sculptor

**Η Ελλάδα των τεχνών και της επιστήμης**  
 Η Ελλάδα αποτελεί έναν χώρο πολιτισμού, τέχνης και επιστημών. Η μακραίωνη συμβολή της στο παγκόσμιο γίνεσθαι, σε συνδυασμό με το μοναδικό φυσικό κάλλος και τις **άρτιες υποδομές**, την καθιστούν ιδανικό τόπο διεξαγωγής συνεδρίων. Τα τελευταία χρόνια, η ελληνική επικράτεια υποδέχεται όλο και συχνότερα ανθρώπους των γραμμάτων, των επιστημών και των τεχνών, οι οποίοι συμμετέχουν σε **συμπόσια, συνέδρια** και **εκθέσεις**. Ο **Διεθνής Αερολιμένας Αθηνών** «Ελευθέριος Βενιζέλος», ένα από τα πλέον σύγχρονα αεροδρόμια παγκοσμίως, ο οποίος λειτουργεί από το 2001, έδωσε μεγάλη ώθηση στη διοργάνωση διεθνών συνεδρίων.

Ο συνεδριακός τουρισμός είναι άκρως αλληλεπιδραστικός: απαιτεί, βέβαια, ένα υψηλό επίπεδο υπόβαθρο από τη χώρα υποδοχής, ταυτόχρονα όμως συμβάλλει ενεργά στην αναβάθμιση της συνολικής ποιότητας μιας περιοχής. Είναι λογικό, ένας χώρος ο οποίος προτιμάται για τη διεξαγωγή συνεδρίων, να μετέχει προνομιακά στο πολιτιστικό «προϊόν», μιας και δίνει τη δυνατότητα σε κοινό, κατοίκους και επισκέπτες, να έρθουν σε επαφή με τα ανθρώπινα επιτεύγματα και τις καινοτομίες.

**Η Ελλάδα των προσωκρατικών φιλοσόφων, των μεγάλων ποιητών, του Φειδία και του**

# ...identifiziert die Sprache...



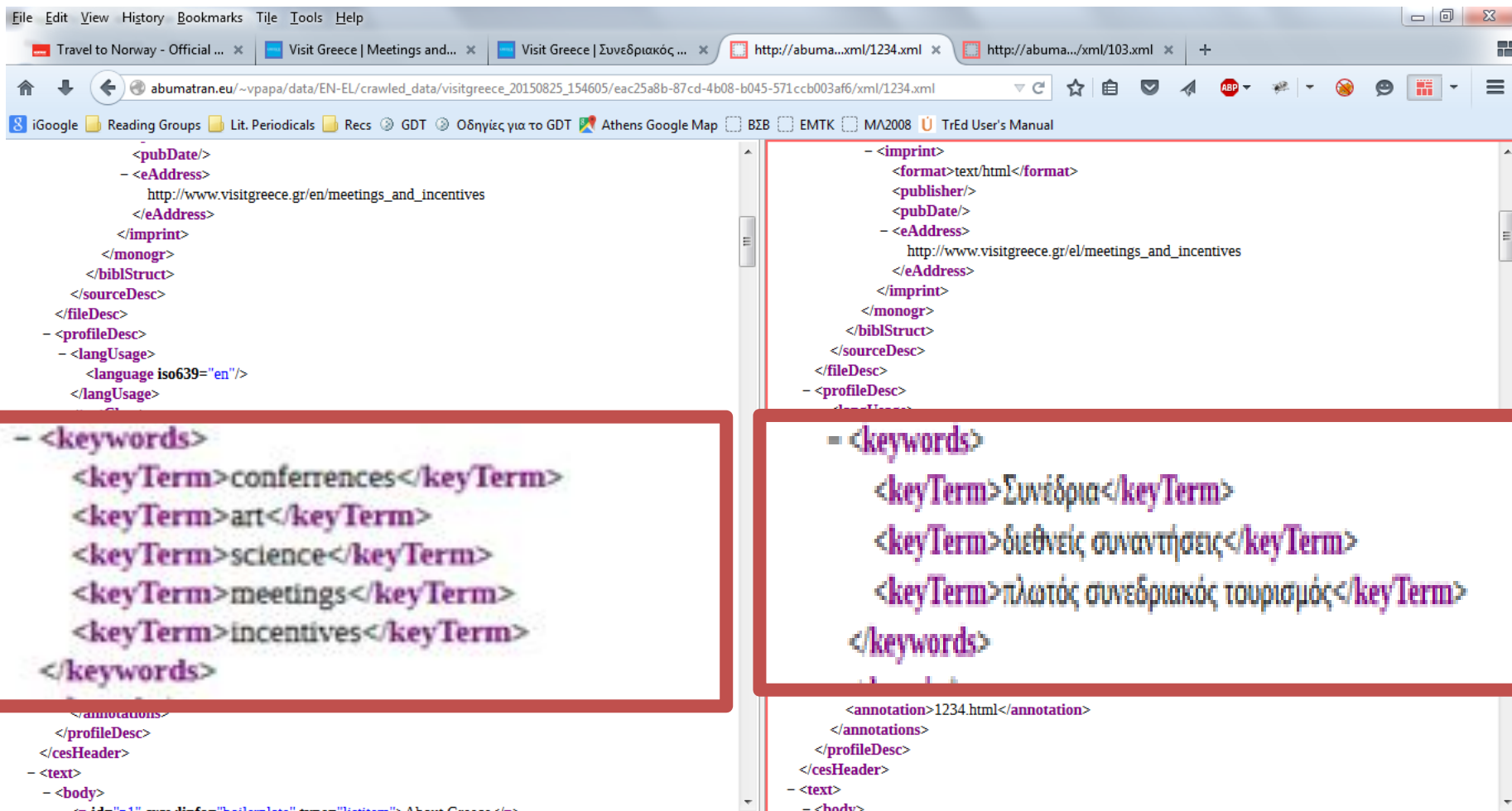
The screenshot displays a web browser window with two XML documents side-by-side. The browser's address bar shows the URL `http://abuma...xml/1234.xml`. The left pane shows the XML code for the English version, with a red box highlighting the following snippet:

```
</fileDesc>
- <profileDesc>
  - <langUsage>
    <language iso639="en"/>
  </langUsage>
```

The right pane shows the XML code for the Greek version, with a red box highlighting the following snippet:

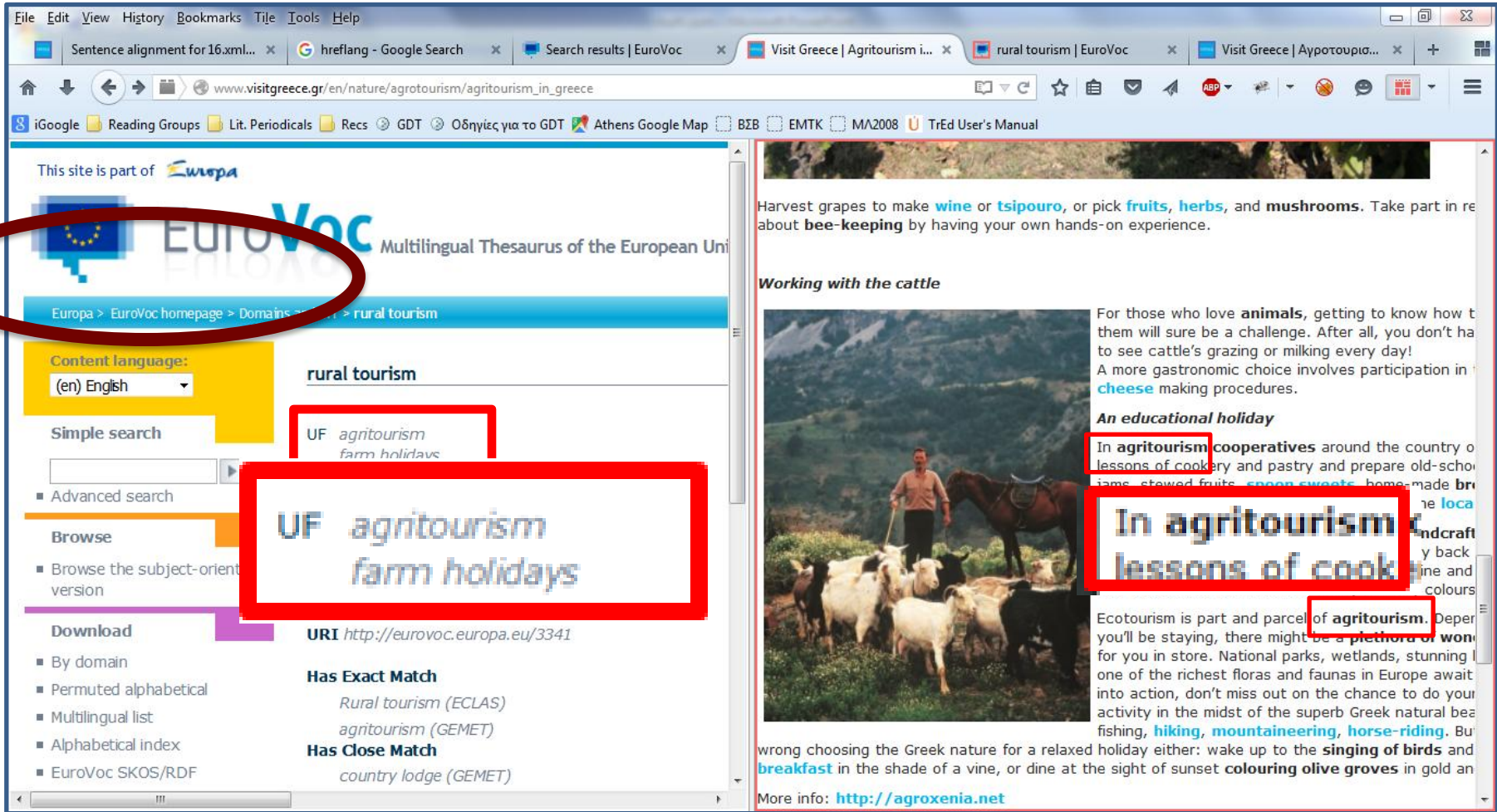
```
- <profileDesc>
  - <langUsage>
    <language iso639="el"/>
  </langUsage>
```

```
Source of: http://www.visitnorway.com/uk/getting-here-and-around/ - Mozilla Firefox
File Edit View Help
49
50 <link rel="alternate" hreflang="da" href="http://www.visitnorway.com/dk/transportmuligheder/" />
51
52 <link rel="alternate" hreflang="es" href="http://www.visitnorway.com/es/como-llegar-y-como-
moverse/" />
53
54 <link rel="alternate" hreflang="fr" href="http://www.visitnorway.com/fr/venir-et-se-deplacer-
en-norvege/" />
55
56 <link rel="alternate" hreflang="nl" href="http://www.visitnorway.com/nl/reizen-naar-en-
in-noorwegen/" />
57
58 <link rel="alternate" hreflang="no" href="http://www.visitnorway.com/no/transport/" />
59
60 <link rel="alternate" hreflang="sv" href="http://www.visitnorway.com/se/transport/" />
61
62 <link rel="alternate" hreflang="it" href="http://www.visitnorway.com/it/arrivare-e-muoversi/" />
63
64 <link rel="alternate" hreflang="ru" href="http://www.visitnorway.com/ru/getting-here-and-
around/" />
65
66 <link rel="alternate" hreflang="en-US" href="http://www.visitnorway.com/us/getting-here-and-
around/" />
67
68 <link rel="alternate" hreflang="pl" href="http://www.visitnorway.com/pl/transport-w-norwegii/"
/>
69
70 <link rel="alternate" hreflang="zh-CN" href="http://www.visitnorway.com/cn/getting-here-and-
around/" />
71
72 <link rel="alternate" hreflang="pt-BR" href="http://www.visitnorway.com/br/como-chegar-
e-mover-se/" />
73
74
75
76 <link rel="stylesheet" href="//d1fy7ceogli51g.cloudfront.net/bundles/styles/visit-
```



```

<pubDate/>
- <eAddress>
  http://www.visitgreece.gr/en/meetings_and_incentives
</eAddress>
</imprint>
</monogr>
</biblStruct>
</sourceDesc>
</fileDesc>
- <profileDesc>
  - <langUsage>
    <language iso639="en"/>
  </langUsage>
- <keywords>
  <keyTerm>conferences</keyTerm>
  <keyTerm>art</keyTerm>
  <keyTerm>science</keyTerm>
  <keyTerm>meetings</keyTerm>
  <keyTerm>incentives</keyTerm>
</keywords>
</annotations>
</profileDesc>
</cesHeader>
- <text>
  - <body>
    <div class="text" data-bbox="20 217 970 897">
      - <imprint>
        <format>text/html</format>
        <publisher/>
        <pubDate/>
        - <eAddress>
          http://www.visitgreece.gr/el/meetings_and_incentives
        </eAddress>
      </imprint>
      </monogr>
      </biblStruct>
      </sourceDesc>
      </fileDesc>
      - <profileDesc>
        <annotation>1234.html</annotation>
      </annotations>
    </profileDesc>
  </cesHeader>
</text>
</body>
  
```

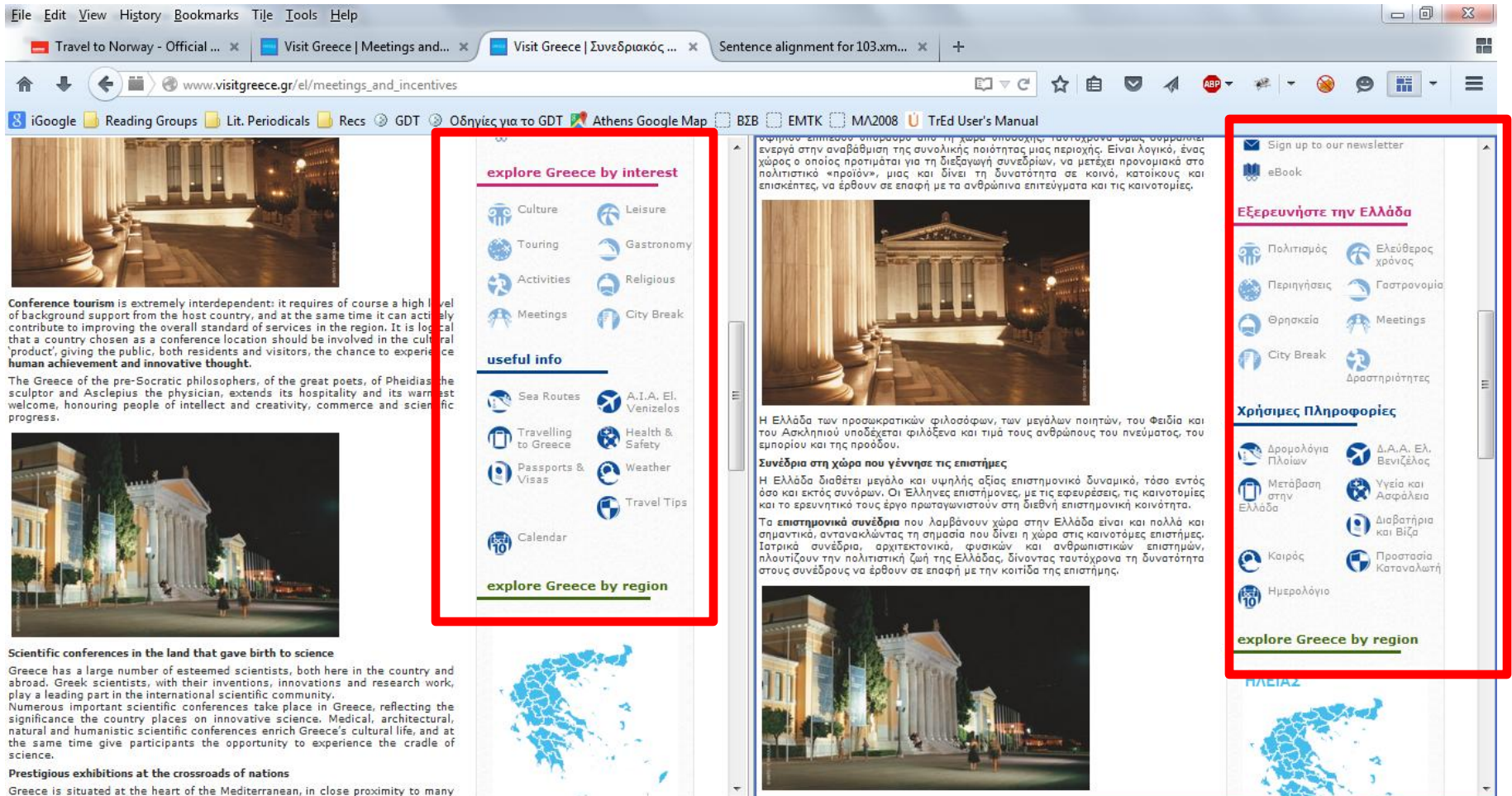



The screenshot shows a web browser window displaying the EuroVoc website. The browser's address bar shows the URL [www.visitgreece.gr/en/nature/agrotourism/agritourism\\_in\\_greece](http://www.visitgreece.gr/en/nature/agrotourism/agritourism_in_greece). The page content includes a search bar with the text "rural tourism" and a list of search results. The results are categorized into "UF" (Upper Field) and "Has Exact Match" and "Has Close Match".

Annotations on the screenshot include:

- A red oval around the EuroVoc logo and navigation menu.
- A red box around the search results for "UF agritourism" and "farm holidays".
- A red box around the text "In agritourism cooperatives" in the article snippet.
- A red box around the text "In agritourism lessons of cookery" in the article snippet.
- A red box around the text "of agritourism" in the article snippet.

The article snippet on the right side of the page discusses agritourism, mentioning activities like grape harvesting, bee-keeping, and working with cattle. It also mentions that agritourism is part of ecotourism and provides information on how to find agritourism cooperatives and lessons of cookery.

File Edit View History Bookmarks Title Tools Help

Travel to Norway - Official ... Visit Greece | Meetings and... Visit Greece | Συνεδριακός ... Sentence alignment for 103.xm...

www.visitgreece.gr/el/meetings\_and\_incentives

iGoogle Reading Groups Lit. Periodicals Recs GDT Οδηγίες για το GDT Athens Google Map BZB EMTK ΜΑ2008 TrEd User's Manual

**explore Greece by interest**

- Culture
- Leisure
- Touring
- Gastronomy
- Activities
- Religious
- Meetings
- City Break

**useful info**

- Sea Routes
- A.I.A. El. Venizelos
- Travelling to Greece
- Health & Safety
- Passports & Visas
- Weather
- Travel Tips
- Calendar

**explore Greece by region**

**ΕΞερευνήστε την Ελλάδα**

- Πολιτισμός
- Ελεύθερος χρόνος
- Περιηγήσεις
- Γαστρονομία
- Θρακεία
- Meetings
- City Break
- Δραστηριότητες

**Χρήσιμες Πληροφορίες**

- Δρομολόγια Πλοίων
- Δ.Α.Α. Ελ. Βενιζέλος
- Μετάβαση στην Ελλάδα
- Υγεία και Ασφάλεια
- Διαβατήρια και Βίζα
- Καιρός
- Προστασία Καταναλωτή
- Ημερολόγιο

**explore Greece by region**

Sign up to our newsletter

eBook

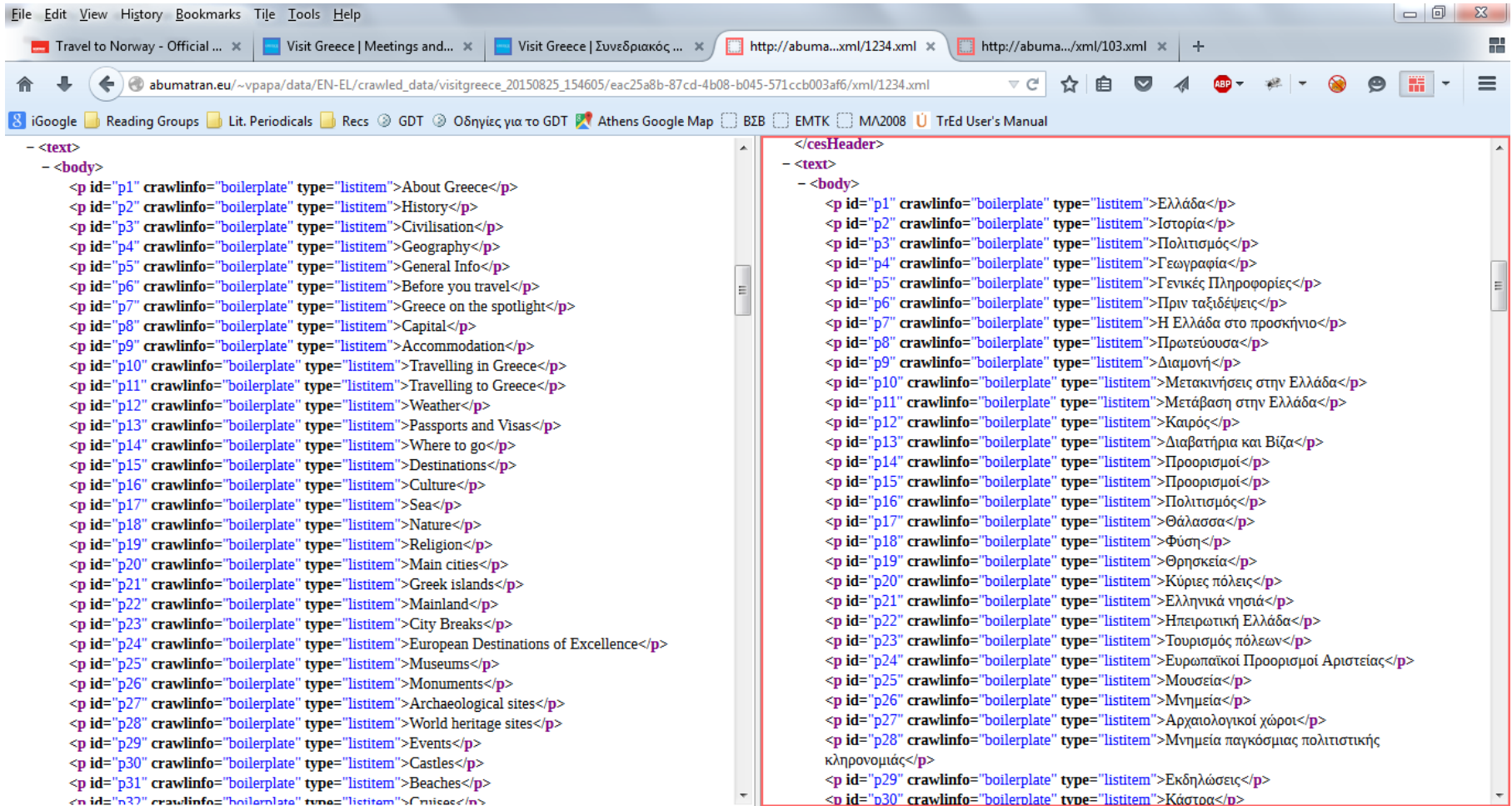
οφίτιος επίκεντρο ομορφιάς από τη χώρα ομορφιάς ταυτόχρονα όμως συμβάλλει ενεργά στην αναβάθμιση της συνολικής ποιότητας μιας περιοχής. Είναι λογικό, ένας χώρος ο οποίος προτιμάται για τη διεξαγωγή συνεδρίων, να μετέχει προνομιακά στο πολιτιστικό «προϊόν», μιας και δίνει τη δυνατότητα σε κοινό, κατοίκους και επισκέπτες, να έρθουν σε επαφή με τα ανθρώπινα επιτεύγματα και τις καινοτομίες.

Η Ελλάδα των προσφρακτικών φιλοσόφων, των μεγάλων ποιητών, του Φειδία και του Ασκληπιού υποδέχεται φιλόξενα και τιμά τους ανθρώπους του πνεύματος, του εμπορίου και της προόδου.

**Συνέδρια στη χώρα που γέννησε τις επιστήμες**

Η Ελλάδα διαθέτει μεγάλο και υψηλής αξίας επιστημονικό δυναμικό, τόσο εντός όσο και εκτός συνόρων. Οι Έλληνες επιστήμονες, με τις εφευρέσεις, τις καινοτομίες και το ερευνητικό τους έργο πρωταγωνιστούν στη διεθνή επιστημονική κοινότητα.

Τα **επιστημονικά συνέδρια** που λαμβάνουν χώρα στην Ελλάδα είναι και πολλά και σημαντικά, αντανακλώντας τη σημασία που δίνει η χώρα στις καινοτόμες επιστήμες. Ιστορικά συνέδρια, αρχιτεκτονικά, φυσικά και ανθρωπιστικών επιστημών, πλουτίζουν την πολιτιστική ζωή της Ελλάδας, δίνοντας ταυτόχρονα τη δυνατότητα στους συνέδρους να έρθουν σε επαφή με την κοίτη της επιστήμης.



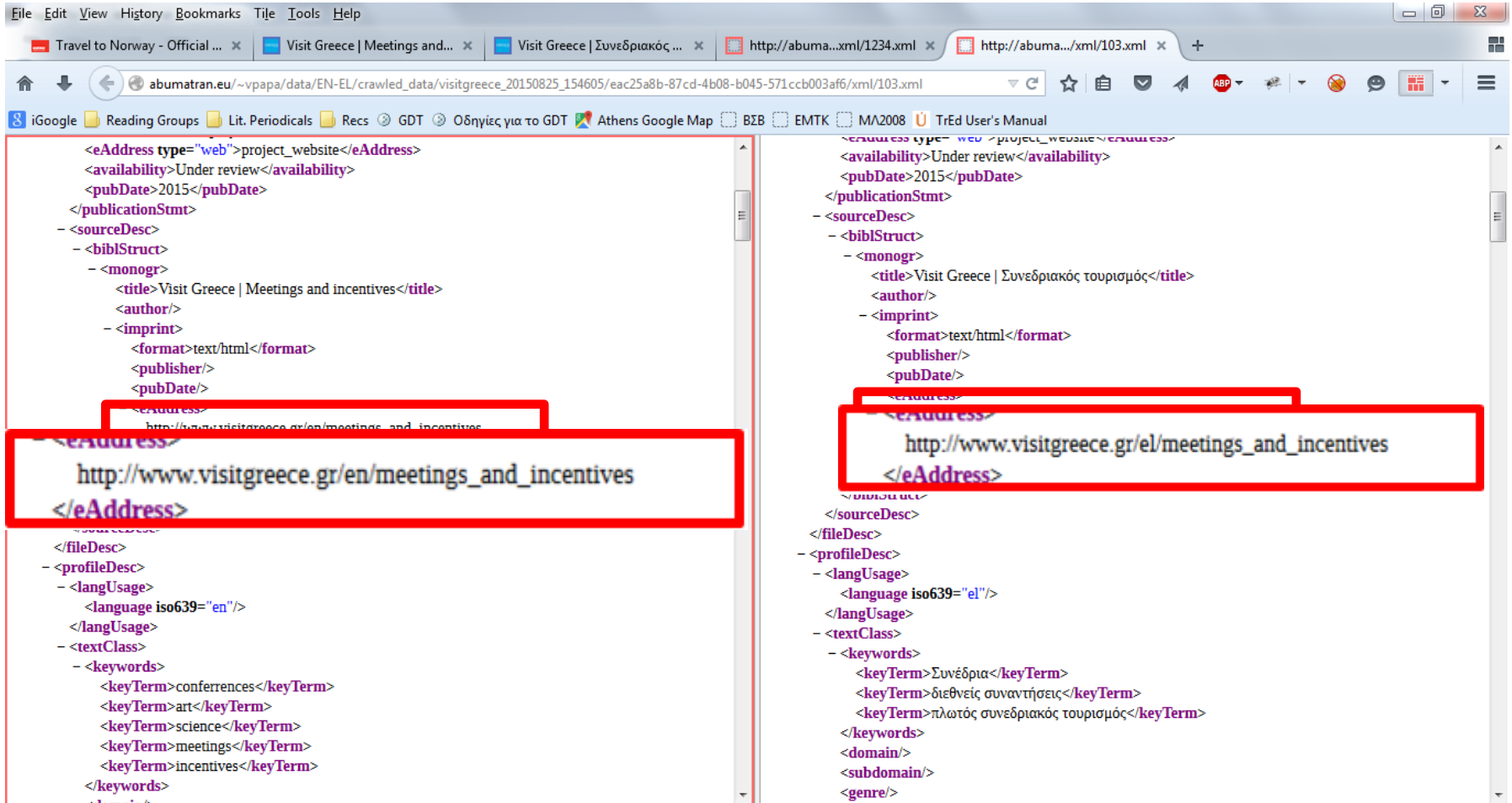
```

- <text>
- <body>
  <p id="p1" crawlinfo="boilerplate" type="listitem">About Greece</p>
  <p id="p2" crawlinfo="boilerplate" type="listitem">History</p>
  <p id="p3" crawlinfo="boilerplate" type="listitem">Civilisation</p>
  <p id="p4" crawlinfo="boilerplate" type="listitem">Geography</p>
  <p id="p5" crawlinfo="boilerplate" type="listitem">General Info</p>
  <p id="p6" crawlinfo="boilerplate" type="listitem">Before you travel</p>
  <p id="p7" crawlinfo="boilerplate" type="listitem">Greece on the spotlight</p>
  <p id="p8" crawlinfo="boilerplate" type="listitem">Capital</p>
  <p id="p9" crawlinfo="boilerplate" type="listitem">Accommodation</p>
  <p id="p10" crawlinfo="boilerplate" type="listitem">Travelling in Greece</p>
  <p id="p11" crawlinfo="boilerplate" type="listitem">Travelling to Greece</p>
  <p id="p12" crawlinfo="boilerplate" type="listitem">Weather</p>
  <p id="p13" crawlinfo="boilerplate" type="listitem">Passports and Visas</p>
  <p id="p14" crawlinfo="boilerplate" type="listitem">Where to go</p>
  <p id="p15" crawlinfo="boilerplate" type="listitem">Destinations</p>
  <p id="p16" crawlinfo="boilerplate" type="listitem">Culture</p>
  <p id="p17" crawlinfo="boilerplate" type="listitem">Sea</p>
  <p id="p18" crawlinfo="boilerplate" type="listitem">Nature</p>
  <p id="p19" crawlinfo="boilerplate" type="listitem">Religion</p>
  <p id="p20" crawlinfo="boilerplate" type="listitem">Main cities</p>
  <p id="p21" crawlinfo="boilerplate" type="listitem">Greek islands</p>
  <p id="p22" crawlinfo="boilerplate" type="listitem">Mainland</p>
  <p id="p23" crawlinfo="boilerplate" type="listitem">City Breaks</p>
  <p id="p24" crawlinfo="boilerplate" type="listitem">European Destinations of Excellence</p>
  <p id="p25" crawlinfo="boilerplate" type="listitem">Museums</p>
  <p id="p26" crawlinfo="boilerplate" type="listitem">Monuments</p>
  <p id="p27" crawlinfo="boilerplate" type="listitem">Archaeological sites</p>
  <p id="p28" crawlinfo="boilerplate" type="listitem">World heritage sites</p>
  <p id="p29" crawlinfo="boilerplate" type="listitem">Events</p>
  <p id="p30" crawlinfo="boilerplate" type="listitem">Castles</p>
  <p id="p31" crawlinfo="boilerplate" type="listitem">Beaches</p>
  <p id="p32" crawlinfo="boilerplate" type="listitem">Cruises</p>
</body>
</cesHeader>
- <text>
- <body>
  <p id="p1" crawlinfo="boilerplate" type="listitem">Ελλάδα</p>
  <p id="p2" crawlinfo="boilerplate" type="listitem">Ιστορία</p>
  <p id="p3" crawlinfo="boilerplate" type="listitem">Πολιτισμός</p>
  <p id="p4" crawlinfo="boilerplate" type="listitem">Γεωγραφία</p>
  <p id="p5" crawlinfo="boilerplate" type="listitem">Γενικές Πληροφορίες</p>
  <p id="p6" crawlinfo="boilerplate" type="listitem">Πριν ταξιδέψεις</p>
  <p id="p7" crawlinfo="boilerplate" type="listitem">Η Ελλάδα στο προσκήνιο</p>
  <p id="p8" crawlinfo="boilerplate" type="listitem">Πρωτεύουσα</p>
  <p id="p9" crawlinfo="boilerplate" type="listitem">Διαμονή</p>
  <p id="p10" crawlinfo="boilerplate" type="listitem">Μετακινήσεις στην Ελλάδα</p>
  <p id="p11" crawlinfo="boilerplate" type="listitem">Μετάβαση στην Ελλάδα</p>
  <p id="p12" crawlinfo="boilerplate" type="listitem">Καιρός</p>
  <p id="p13" crawlinfo="boilerplate" type="listitem">Διαβατήρια και Βίζα</p>
  <p id="p14" crawlinfo="boilerplate" type="listitem">Προορισμοί</p>
  <p id="p15" crawlinfo="boilerplate" type="listitem">Προορισμοί</p>
  <p id="p16" crawlinfo="boilerplate" type="listitem">Πολιτισμός</p>
  <p id="p17" crawlinfo="boilerplate" type="listitem">Θάλασσα</p>
  <p id="p18" crawlinfo="boilerplate" type="listitem">Φύση</p>
  <p id="p19" crawlinfo="boilerplate" type="listitem">Θρησκεία</p>
  <p id="p20" crawlinfo="boilerplate" type="listitem">Κύριες πόλεις</p>
  <p id="p21" crawlinfo="boilerplate" type="listitem">Ελληνικά νησιά</p>
  <p id="p22" crawlinfo="boilerplate" type="listitem">Ηπειρωτική Ελλάδα</p>
  <p id="p23" crawlinfo="boilerplate" type="listitem">Τουρισμός πόλεων</p>
  <p id="p24" crawlinfo="boilerplate" type="listitem">Ευρωπαϊκοί Προορισμοί Αριστείας</p>
  <p id="p25" crawlinfo="boilerplate" type="listitem">Μουσεία</p>
  <p id="p26" crawlinfo="boilerplate" type="listitem">Μνημεία</p>
  <p id="p27" crawlinfo="boilerplate" type="listitem">Αρχαιολογικοί χώροι</p>
  <p id="p28" crawlinfo="boilerplate" type="listitem">Μνημεία παγκόσμιας πολιτιστικής κληρονομιάς</p>
  <p id="p29" crawlinfo="boilerplate" type="listitem">Εκδηλώσεις</p>
  <p id="p30" crawlinfo="boilerplate" type="listitem">Κάστρα</p>
  <p id="p31" crawlinfo="boilerplate" type="listitem">Κρουαζιέρες</p>
  <p id="p32" crawlinfo="boilerplate" type="listitem">Κρουαζιέρες</p>
</body>

```



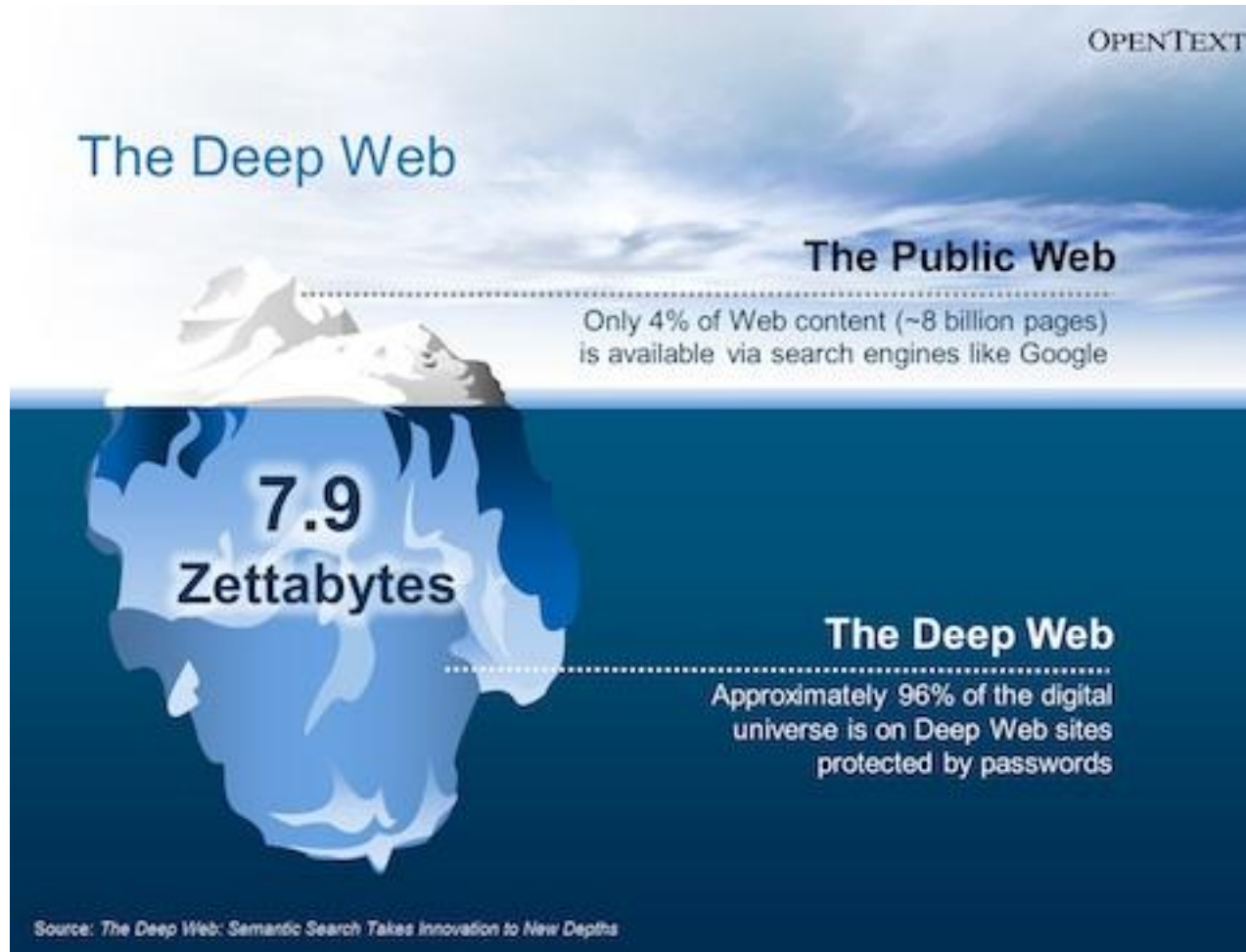
# ... nutz HTML Struktur und URL Ähnlichkeiten zur Erkennung von Dokumentenpaaren...



The screenshot shows a web browser with two XML documents open side-by-side. The left document is from `abumatran.eu/~vpapa/data/EN-EL/crawled_data/visitgreece_20150825_154605/eac25a8b-87cd-4b08-b045-571ccb003af6/xml/103.xml` and the right document is from `http://abuma.../xml/103.xml`. Both documents contain XML metadata for a document titled "Visit Greece | Meetings and incentives" (left) and "Visit Greece | Συνεδριακός τουρισμός" (right). The URL `http://www.visitgreece.gr/en/meetings_and_incentives` is highlighted in red in both documents, demonstrating a similarity in the HTML structure used to identify document pairs.

| #  | en   | de  |
|----|--|---|
| 1  | Reliable partnership based on trust  | Vertrauensvolle und verlässliche Partnerschaft  |
| 2  | "The stage has been set for further intensive cooperation," declared Chancellor <b>Angela Merkel</b> after her meeting with Colombian President <b>Juan Manuel Santos</b> .  | "Die Weichen für eine weitere intensive Zusammenarbeit sind gestellt." Das hat Bundeskanzlerin <b>Merkel</b> nach dem Treffen mit dem kolumbianischen Präsidenten <b>Santos</b> erklärt.  |
| 3  | The Chancellor intends to support the peace process in Colombia, partly through cooperation in the fields of research, education and climate change mitigation.  | Den Friedensprozess in Kolumbien will Merkel unter anderem durch weitere Kooperationen in den Bereichen Forschung, Bildung und Klimaschutz unterstützen.  |
| 4  | Colombia's President seeks support in the peace process with FARC rebels Photo:  | Kolumbiens Präsident wirbt für Unterstützung beim Friedensprozess mit den FARC-Rebellen. Foto:  |
| 5  | Bundesregierung/Denzel   | Bundesregierung/Denzel  |
| 6  | "We have a cordial and reliable partnership based on trust," declared Chancellor <b>Angela Merkel</b> after her meeting with Colombian President <b>Juan Manuel Santos</b> .   | "Wir sind in einer vertrauensvollen, freundschaftlichen und verlässlichen Partnerschaft", erklärte Bundeskanzlerin <b>Angela Merkel</b> nach dem Treffen mit dem kolumbianischen Präsidenten <b>Juan Manuel Santos</b> .  |
| 7  | Their talks focused on Colombia's peace process, bilateral relations and economic and regional issues.   | Im Mittelpunkt des Gesprächs standen der kolumbianische Friedensprozess, die bilateralen Beziehungen sowie wirtschaftliche und regionalpolitische Themen.   |
| 8  | Fostering the peace process  | Friedensprozess fördern   |
| 9  | "The situation today in Colombia is marked by the courageous peace process initiated by the President, which is currently in a crucial phase," said the Chancellor following talks with Colombia's President Juan Manuel Santos. | "Die aktuelle Situation in Kolumbien ist dadurch gekennzeichnet, dass der Präsident einen mutigen Friedensprozess initiiert hat, der im Augenblick in einer entscheidenden Phase ist", erläuterte die Kanzlerin nach dem Treffen mit dem kolumbischen Präsidenten Juan Manuel Santos. |
| 10 | She reported that she had pledged Germany's full support in this process.  | Hierfür habe sie die volle deutsche Unterstützung zugesagt, sagte Merkel.   |
| 11 | In Colombia a conflict has been smouldering for decades between right-wing paramilitary groups, left-wing guerrillas and the Colombian army.   | In Kolumbien schwelt seit Jahrzehnten ein Konflikt zwischen rechtsgerichteten Paramilitärs, linksgerichteten Guerillatruppen und der kolumbianischen Armee.   |

- Wir können den nach außen sichtbaren Teil der Daten finden, aber es gibt noch viel mehr in Ihren Organisationen
- Helfen Sie uns diese Daten zu bekommen und zu nutzen!
- Dieser Prozess kann mit Ihrer Hilfe zu einer automatisierten Pipeline zur Herstellung von parallelen Sprachressourcen werden (sammeln Sie Ihre Dokumente, Berichte, Dateien, etc.)

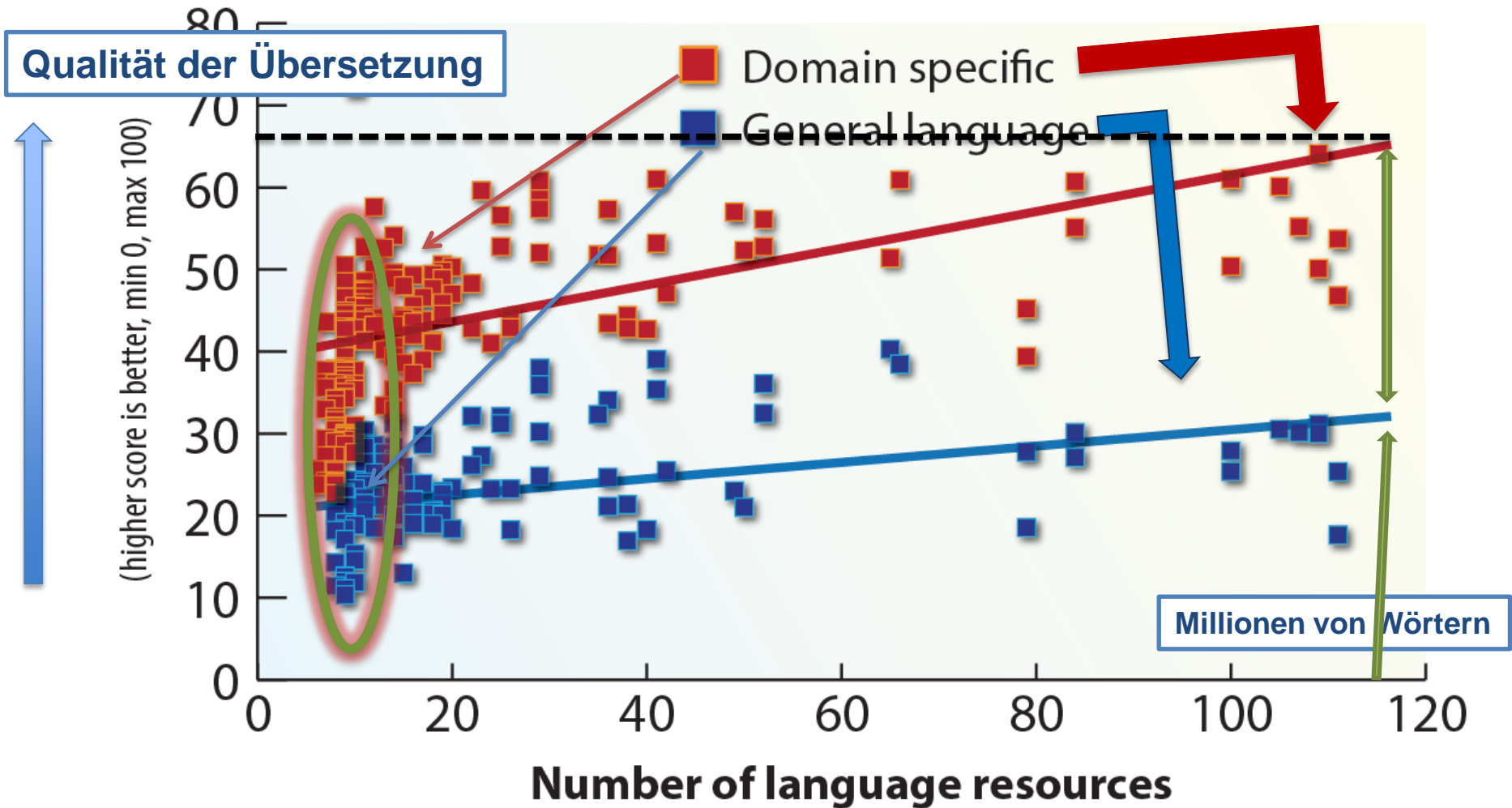






- Solche Dokumente gibt es in:
  - Dokumentationszentren (übersetzte Berichte, Broschüren, Vorträge, (interne) Webseiten, etc.)
  - Bei den Language Service Providers (LSP), welche Übersetzungsaufträge erhalten
- Helfen Sie uns diese Quellen zu identifizieren und mit diesen zu kooperieren

# Impact of number of language resources on Statistical MT quality



- Verwertung von bereits vorhandenen Daten
- Effiziente Identifikation von Quellen ist
- Der Wert der vorhandenen Ressourcen darf nicht unterschätzt werden
- Der Startpunkt der Suche ist oft das Schwierigste

**Zusammenarbeit in der  
Auffindung der Quellen ist essentiell**