

SULLA CORREZIONE ORTOGRAFICA AUTOMATICA

Un'applicazione concreta partendo da Wikipedia

Luca Chiodini

Esame di Stato 2015

ITIS Paleocapa

The image shows a screenshot of a Google search interface. At the top left is the Google logo. The search bar contains the text "prvisioni mteoa". To the right of the search bar are a microphone icon and a blue search button with a magnifying glass icon. Below the search bar, there are navigation tabs: "Web" (highlighted with a red underline), "Notizie", "Shopping", "Immagini", "Maps", "Altro" (with a downward arrow), and "Strumenti di ricerca". Below the tabs, it says "Circa 6.530.000 risultati (0,51 secondi)". Underneath, there is a heading "Risultati relativi a *previsioni meteo*" and a sub-heading "Cerca invece [prvisioni mteoa](#)".

COSA C'È DIETRO LE QUINTE?

amre

amre

Amore? Mare? Amare?

amre

Amore? Mare? Amare?

$$\arg \max_{c \in D} \underbrace{P(w|c)}_{\text{error}} \cdot \underbrace{P(c)}_{\text{language}}$$

nuotano nel amre

nuotano nel amre

Amore? Mare? Amare?

l'amre è cieco

Amore? Mare? Amare?

l'amre è cieco

Amore? Mare? Amare?

$$\arg \max_{c \in D} \underbrace{P(w|c_i)}_{\text{error}} \cdot \underbrace{P(c_{i-1}|c_i) \cdot P(c_i) \cdot P(c_i|c_{i+1}))}_{\text{language}}$$

- Wikipedia in italiano come sorgente di dati

- Wikipedia in italiano come sorgente di dati
- Analizzare **tutta** Wikipedia **non è banale!**

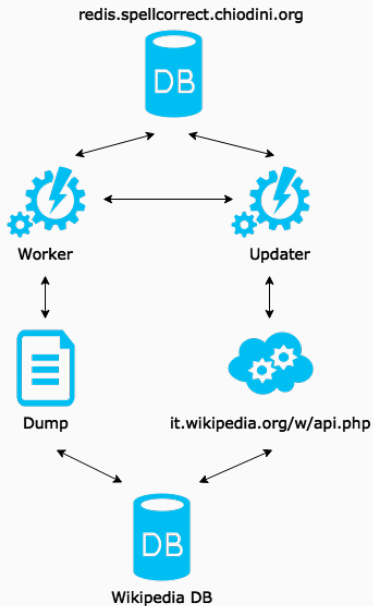
- Wikipedia in italiano come sorgente di dati
- Analizzare **tutta** Wikipedia **non è banale!**
- ~ 10 GB di pagine (solo testo)

- Wikipedia in italiano come sorgente di dati
- Analizzare **tutta** Wikipedia **non è banale!**
- ~ 10 GB di pagine (solo testo)
- Elaborazione su Amazon EC2 (8 CPU, 15 GiB RAM)

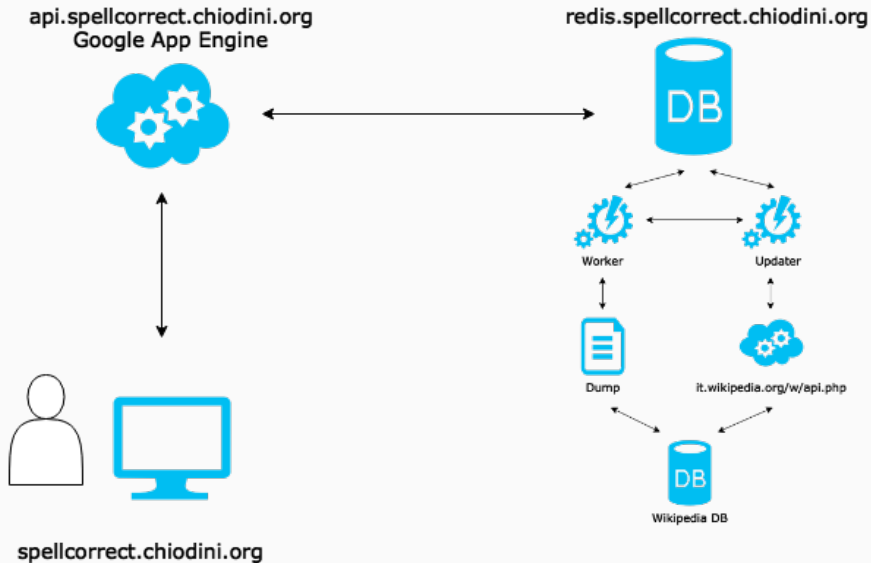
- Wikipedia in italiano come sorgente di dati
- Analizzare **tutta** Wikipedia **non è banale!**
- ~ 10 GB di pagine (solo testo)
- Elaborazione su Amazon EC2 (8 CPU, 15 GiB RAM)
- ~ 311 000 000 parole (~ 2 600 000 diverse)

- Wikipedia in italiano come sorgente di dati
- Analizzare **tutta** Wikipedia **non è banale!**
- ~ 10 GB di pagine (solo testo)
- Elaborazione su Amazon EC2 (8 CPU, 15 GiB RAM)
- ~ 311 000 000 parole (~ 2 600 000 diverse)
- ~ 41 000 000 coppie di parole diverse

ARCHITETTURA DEL SISTEMA



ARCHITETTURA DEL SISTEMA



Tutto ciò che è bello e nobile è il risultato della ragione e di calcoli.

Baudelaire