#### **GABRIELE ROMANATO**

Menu

#### Come estrarre il testo da un documento PDF con Go

L'estrazione del testo da un documento PDF è una necessità comune in molti scenari, come il data mining, l'automazione aziendale o l'analisi dei dati. In questo articolo esploreremo come realizzare un semplice programma in Go per estrarre il testo da un file PDF. Go è un linguaggio di programmazione moderno e potente, particolarmente adatto per lo sviluppo di applicazioni ad alte prestazioni. Per completare questa operazione useremo una libreria open-source di Go chiamata pdfcpu, che consente di gestire i file PDF, incluso l'estrazione del testo.

Prima di tutto, è necessario installare il package pdfcpu nel progetto Go corrente. Si può fare tramite il comando:

```
go get github.com/pdfcpu/pdfcpu
```

Ecco un semplice esempio di come estrarre il testo da un PDF utilizzando pdfcpu nel file main. go del progetto:

```
package main
import (
   "fmt"
   "log"
   "os"
   "github.com/pdfcpu/pdfcpu/pkg/api"
    "github.com/pdfcpu/pdfcpu/pkg/pdfcpu"
)
func main() {
   // Verifica se il file PDF è stato passato come argomento
   if len(os.Args) < 2 {
        log.Fatalf("Utilizzo: %s <file.pdf>\n", os.Args[0])
   }
   pdfFile := os.Args[1]
   // Configura PDFCPU per estrarre il testo
   conf := pdfcpu.NewDefaultConfiguration()
   // Estrazione del testo dal PDF
   out, err := api.ExtractTextFile(pdfFile, conf)
   if err != nil {
       log.Fatalf("Errore nell'estrazione del testo: %v\n", err)
   }
   // Stampa il testo estratto
   fmt.Println(string(out))
}
```

Spiegazione del codice:

- Importazione dei pacchetti: Il codice utilizza i pacchetti fmt, log e os forniti da Go, oltre a pdfcpu per manipolare i file PDF.
- Verifica degli argomenti: Il programma verifica se è stato passato un file PDF come argomento. Se non lo è, viene mostrato un messaggio di errore.
- Estrazione del testo: La funzione api.ExtractTextFile viene utilizzata per estrarre il testo dal file PDF. La funzione prende due parametri: il percorso del file PDF e una configurazione di default.
- Stampa del testo estratto: Il testo estratto viene convertito in una stringa e stampato sulla console.

#### Conclusione

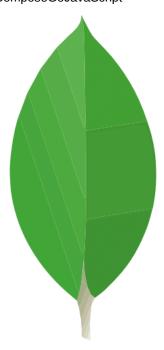
L'estrazione del testo da un PDF in Go è un compito relativamente semplice grazie alla libreria pdfcpu. Con poche righe di codice, puoi estrarre efficacemente il testo da file PDF e utilizzarlo per ulteriori elaborazioni. Questo strumento può essere molto utile per una varietà di applicazioni, specialmente se integrato in pipeline di elaborazione automatizzata o in sistemi di gestione documentale.

## **Applicazioni Correlate**



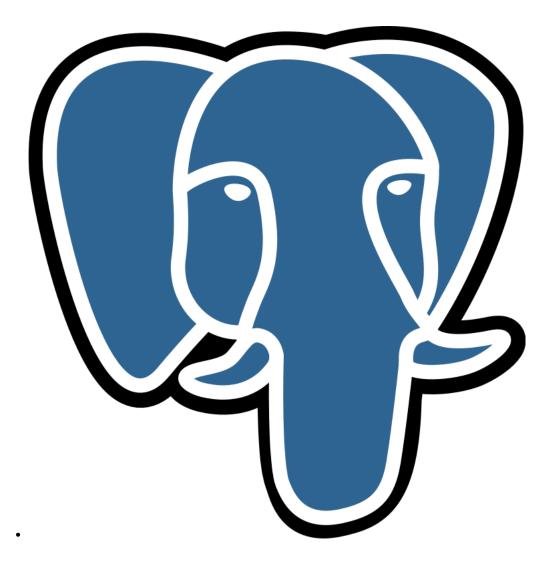
### **Go Placeholder Image**

Applicazione in Go per la creazione di immagini segnaposto. DockerDocker ComposeGoJavaScript



## Go MongoDB App

Applicazione basata su MongoDB ed implementata in Go con il driver ufficiale. DockerDocker ComposeGoMongoDB



# Go PostgreSQL App

Applicazione basata su PostgreSQL e sviluppata in Go. DockerDocker ComposeGoPostgreSQL