

Opgave: PDF Downloader

Afleveringsformat:

Opgaven skal afleveres i form af et link til et GitHub repository hvor I har lavet projektet. Start med at lave en kravspecifikation ud fra skabelonen der er uploadet i mappen for softwareudvikling på Teams.

Derudover skal der mindst være enten et UML-klasse diagram eller et UML-sekvens diagram der viser systemets opbygning. Yderligere skal der være kommentarer i koden for at gøre den lettere at forstå. Din endelige løsning skal inkludere en Readme-guide.

Historie:

I din rolle som IT-konsulent hos Specialisterne har du modtaget en opgave fra en kunde, der har et ældre Python-script, som ikke længere fungerer korrekt. Kunden har brug for et pålideligt og effektivt program til at downloade PDF-rapporter fra en liste med dynamiske URL'er. Der er behov for at håndtere alternative URL'er, hvor det første link ikke virker, og at registrere downloadstatus for hver rapport.

Beskrivelse fra Kunden:

"Kære Konsulent,

Som aftalt har du hermed listen (inkl. metadata) angående de rapporter vi gerne vil have downloaded ("GRI_2017_2020").

Det er adressen i kolonne AL (Pdf_URL) vi har forsøgt at downloade. Men som jeg forklarede ville det være fedt, at hvis dette link ikke virker, så prøvede programmet linket i kolonne AM. Hvis første link virker, behøver den ikke prøve AM.

Jeg har også vedhæftet vores Python-program, som vi har brugt til at downloade med. Som sagt kører koden, men den er meget ustabil, og det går til tider meget langsomt.

Desuden er der vedhæftet en Excel-fil med rapporter fra 2006-2016 ("Metadata2006-2016"), hvor man i kolonne AT kan se, om en rapport er blevet downloadet eller ej (som vi også gerne vil have information om med de nye rapporter). Vi behøver ikke en specificering af, hvorfor rapporten ikke er hentet, blot en variabel der f.eks. antager værdien "Downloadet" eller "Ikke downloadet".

Kravspecifikationen for dette projekt er, at vi ønsker, at I laver et program, der effektivt kan downloade alle de rapporter, der har et virkende link. Programmet skal downloade disse PDF-rapporter og dele dem med os via NAS (eller anden måde). De downloadede rapporter skal navngives efter kolonnen "BRNummer". Derudover ønsker vi en liste over, hvilke rapporter (fra GRI_2017_2020), der er blevet downloadet, og hvilke der ikke er. Den endelige kode afleveres ligeledes til os, da der formegentligt kommer flere rapporter, der skal hentes senere.

Bliver koden lavet i Python, eller er det ikke det rigtige program til denne type opgave?”

OBS:

For ikke at overbelaste Specialisternes netværk - Start med at lave en prototype som kun downloader maks. 10 PDF'er af gangen. Når I arbejder hjemme kan I forsøge med det fulde antal PDF'er og fuld concurrency.