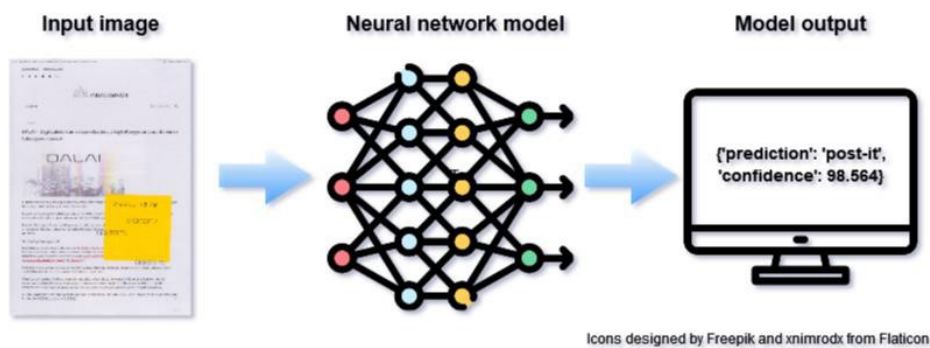


Igenkänning av post-it-lappar

Bakgrund, fördelar och ändamål med komponenten

Vid digitaliseringen av dokument kan det hända att dokument med en eller flera post-it-lappar har överlämnats för skanning. Om post-it-lappens placering påverkar läsbarheten av innehållet, är det bra om sådana fall kan kännas igen och vid behov skannas på nytt. En maskinell igenkänning av post-it-lappar bidrar därmed till att förbättra kvaliteten på digitalt material och minskar behovet av en manuell kvalitetskontroll.

Vad gör komponenten?



Komponenten behandlar de bildfiler som den fått som indata en efter en och tar fram en prognos för bildens innehåll på grundval av de parametrar som den har lärt sig under utbildningen. Om komponenten känner igen en post-it-lapp på bilden med mer än 50 procents sannolikhet, klassificeras bilden som felaktig. Komponentens återställer den förväntade klassen till användaren.

Hur används komponenten?

Komponenten finns tillgänglig både som en del av Arkkiivi-gränssnittet (<http://www.arkkiivi.fi/>) och som en fristående app som tillhandahåller ett programmeringsgränssnitt (API) på DALAI-projektets GitHub-webbplats (<https://github.com/DALAI-project/PostitAPI>). Närmare anvisningar om användningen av komponenten inom olika miljöer finns på den webbplatsen.

Utbildning av komponenten

Bildfiler i .jpg-format har använts för att utbilda komponenten och deras storlek har harmoniserats till 224 x 224 pixlar. Det totala antalet dokument som använts i utbildningen uppgår till cirka 55 000 exemplar, av vilka cirka 4 000 är dokument som innehåller post-it-lappar. Bilder som innehåller post-it-lappar har både gjorts själva och inhämtats från bilder som i samband med Riksarkivets massdigitaliseringsprocess manuellt klassificerats som felaktiga. Dokumenten innehåller ett varierande antal post-it-lappar i olika storlekar och färger, som är placerade på olika sidor av dokumentet och som även kan innehålla text.

Trots att man har försökt samla många olika exempel till utbildningsmaterialet, är materialmängden begränsad och det uppstår också fel i komponentens klassificering. Till exempel kan innehållselement som liknar post-it-lappar, t.ex. fyrkantiga textfälts i färg, leda till en felaktig klassificering.