

**Supplementary File 1: Descriptions of the datasets used in the development, testing and implementation of the occupation application**

<b>Application Development and Testing Datasets</b>			
	<b>Type of document</b>	<b>Document count</b>	<b>No. of Occupation Annotations (manual)</b>
Validation corpus	Personal history	77 (+256 documents used in training)	405
Testing corpus 1: with vs without machine-learning comparison	Personal history + other CRIS documents	200	521
Testing corpus 2: gold-standard annotated documents	Personal history	666	3,429
Testing corpus 4: Unannotated documents	Personal history	200	442
<b>Application Implementation Dataset</b>			
	<b>Type of document</b>	<b>Patient count</b>	<b>No. Of Occupation Extractions (application)</b>
CRIS case register of patient records aged $\geq 16$	Attachments	341,720	21,321,757 (all relations)
	Events		
	Correspondence		
	Discharge Notification Summaries		
	History		

	Mental State Formulations		
	Presenting Circumstances		
	Risk Events		
	Social Situation		
	Ward Progress Notes		

*Table 1: Descriptions of the datasets used in the development, testing and implementation of the occupation application*