

Table 1

<b>Collection</b>	<b>Original documents</b>	<b>Scanned images</b>	<b>Word tokens</b>
<b>British House of Commons Parliamentary Papers (ProQuest)</b>	<b>118,526</b>	<b>6,448,739</b>	<b>4,672,737,402</b>
<b>Early Canada Online (Canadiana.org)</b>	<b>83,016</b>	<b>3,938,758</b>	<b>2,066,746,780</b>
<b>Kew Gardens Directors' Correspondence</b>	<b>24,765</b>		<b>4,550,118</b>
<b>Confidential Prints Collection: Africa, Latin America, Middle East, and North America (Adam Matthew Online)</b>	<b>1,315</b>	<b>140,010</b>	<b>206,389,401</b>
<b>Selected Books and Documents</b>	<b>306</b>		<b>5,123,458</b>
<b>Total</b>	<b>227,928</b>	<b>10,527,507</b>	<b>6,955,547,159</b>

Table 1 Caption:

The Corpus. Scanned images are the pages produced by optical character recognition and generally correspond to pages in the source texts. Word tokens are the strings that are recognized by our text processing software.