

Bunsirou wa tamerawazuni sono yubi wo kuti ni hukumu to, kizuguti wo tuyoku sutta.
文四郎はためらわずにその指を口に含むと、傷口を強く吸った。

BCCWJ sample
OT02_00007

Morphological analysis with
MeCab 0.994 and UniDic 1.3.12

<u>1</u>	<u>2</u>	<u>3</u>	<u>4</u>	<u>5</u>	<u>6</u>
文四郎	は	ためら	わ	ず	に
<u>1</u>		<u>2</u>		<u>3</u>	

characters = 29

one character per Unicode codepoint

tokens = 19
morphemes (short unit words) as defined in
UniDic 1.3.12/Mecab 0.994

<u>7</u>	<u>8</u>	<u>9</u>	<u>10</u>	<u>11</u>	<u>12</u>	<u>13</u>	<u>14</u>	<u>15</u>	<u>16</u>	<u>17</u>	<u>18</u>	<u>19</u>
指	を	口	に	含む	と、	傷口	を	強く	吸	っ	た	。
<u>4</u>		<u>5</u>		<u>6</u>		<u>7</u>		<u>8</u>		<u>9</u>		

chunks = 9

phrasal unit consisting of content word
and functional unit, as defined in
CaboCha; 'bunsetsu'