

# Verwerking van achttiende-eeuws Nederlands met Frog

Erik Tjong Kim Sang  
Meertens Instituut  
erikt(at)xs4all.nl

13 februari 2014

## 1 Inleiding

Frog [1] is een verzameling programma's voor de taalkundige analyse van hedendaagse Nederlandstalige teksten. In het Nederlabproject<sup>1</sup> willen we ook teksten in ouder Nederlands taalkundig analyseren. In dit document gaan we na of Frog achttiende-eeuwse teksten correct kan verwerken. Bij deze test passen we Frog toe op twee oude teksten en een recente tekst.

## 2 Teksten

We hebben de volgende drie teksten geselecteerd voor onze test:

- bladzijde 2 uit *Poetische en Stigtelyke Mengeldigten*, van Adrianus Hardy, uit 1750.  
<http://resolver.kb.nl/resolve?urn=dpo:2002:mpeg21:0012>
- bladzijde 2 uit *Romeinsche Geschiedenissen*, deel 16, van M. Stuart uit het jaar 1800.  
<http://resolver.kb.nl/resolve?urn=dpo:10362:mpeg21:0008>

---

<sup>1</sup>[www.nederlab.nl](http://www.nederlab.nl)

- bladzijde 2 uit *Benoît, of hoe integreert men het gecrêpempapierde kaasdoosje-voorvaderdag in de loft-met-starck...* van Noortje Wiesbauer uit 1993 (als vergelijkingsmateriaal).  
[http://dbnl.nl/tekst/\\_bra004199301\\_01/\\_bra004199301\\_01\\_0064.php](http://dbnl.nl/tekst/_bra004199301_01/_bra004199301_01_0064.php) [p. 161]

De eerste twee teksten verkregen we via de website Early Dutch Books Online<sup>2</sup>. Van deze twee teksten waren de gescande versies beschikbaar (tekst met scanfouten). De derde tekst verkregen we via de website Digitale Bibliotheek voor de Nederlandse Letteren<sup>3</sup>. Van deze tekst was een PDF-versie beschikbaar, waaruit met knippen en plakken een tekstdocument komt worden afgeleid.

### 3 Voorbewerking

Voordat de teksten konden worden aangeboden aan het programma Frog, moesten ze eerst worden omgezet in het FoLiA [2], het XML-formaat wordt gebruikt voor de in- en uitvoerbestanden van Frog. Voor de recente tekst volstond het om een extra programma te schrijven dat paragrafen tussen paragraaftags plaatste en een XML-inleiding en -uitleiding toevoegde.

De twee oudere teksten konden niet op dezelfde manier worden omgezet omdat op de bladzijden meer informatie aanwezig was dan de tekst:

- bladzijdenummers
- hoofdstuktitels (herhaald op elke pagina)
- noten (zowel onder als naast de tekst)
- de eerste lettergreep van de volgende pagina

\*      ABBREVIECEN

175.    vooel des Remercken Suet. — Nieuw-  
1777.    Bije had ten quater gewelt van het op  
wonder sietes verhooren Dierverhoer van twee  
in 1711.    en 1712.    Brevolde  
die oppenre, de Tien die Edele  
van een wij Concellet te woen, wie  
die twee, of 1711.    en 1712.    en 1713.    en 1714.    en 1715.    en 1716.    en 1717.    en 1718.    en 1719.    en 1720.    en 1721.    en 1722.    en 1723.    en 1724.    en 1725.    en 1726.    en 1727.    en 1728.    en 1729.    en 1730.    en 1731.    en 1732.    en 1733.    en 1734.    en 1735.    en 1736.    en 1737.    en 1738.    en 1739.    en 1740.    en 1741.    en 1742.    en 1743.    en 1744.    en 1745.    en 1746.    en 1747.    en 1748.    en 1749.    en 1750.    en 1751.    en 1752.    en 1753.    en 1754.    en 1755.    en 1756.    en 1757.    en 1758.    en 1759.    en 1760.    en 1761.    en 1762.    en 1763.    en 1764.    en 1765.    en 1766.    en 1767.    en 1768.    en 1769.    en 1770.    en 1771.    en 1772.    en 1773.    en 1774.    en 1775.    en 1776.    en 1777.    en 1778.    en 1779.    en 1780.    en 1781.    en 1782.    en 1783.    en 1784.    en 1785.    en 1786.    en 1787.    en 1788.    en 1789.    en 1790.    en 1791.    en 1792.    en 1793.    en 1794.    en 1795.    en 1796.    en 1797.    en 1798.    en 1799.    en 1800.    en 1801.    en 1802.    en 1803.    en 1804.    en 1805.    en 1806.    en 1807.    en 1808.    en 1809.    en 1810.    en 1811.    en 1812.    en 1813.    en 1814.    en 1815.    en 1816.    en 1817.    en 1818.    en 1819.    en 1820.    en 1821.    en 1822.    en 1823.    en 1824.    en 1825.    en 1826.    en 1827.    en 1828.    en 1829.    en 1830.    en 1831.    en 1832.    en 1833.    en 1834.    en 1835.    en 1836.    en 1837.    en 1838.    en 1839.    en 1840.    en 1841.    en 1842.    en 1843.    en 1844.    en 1845.    en 1846.    en 1847.    en 1848.    en 1849.    en 1850.    en 1851.    en 1852.    en 1853.    en 1854.    en 1855.    en 1856.    en 1857.    en 1858.    en 1859.    en 1860.    en 1861.    en 1862.    en 1863.    en 1864.    en 1865.    en 1866.    en 1867.    en 1868.    en 1869.    en 1870.    en 1871.    en 1872.    en 1873.    en 1874.    en 1875.    en 1876.    en 1877.    en 1878.    en 1879.    en 1880.    en 1881.    en 1882.    en 1883.    en 1884.    en 1885.    en 1886.    en 1887.    en 1888.    en 1889.    en 1890.    en 1891.    en 1892.    en 1893.    en 1894.    en 1895.    en 1896.    en 1897.    en 1898.    en 1899.    en 1900.    en 1901.    en 1902.    en 1903.    en 1904.    en 1905.    en 1906.    en 1907.    en 1908.    en 1909.    en 1910.    en 1911.    en 1912.    en 1913.    en 1914.    en 1915.    en 1916.    en 1917.    en 1918.    en 1919.    en 1920.    en 1921.    en 1922.    en 1923.    en 1924.    en 1925.    en 1926.    en 1927.    en 1928.    en 1929.    en 1930.    en 1931.    en 1932.    en 1933.    en 1934.    en 1935.    en 1936.    en 1937.    en 1938.    en 1939.    en 1940.    en 1941.    en 1942.    en 1943.    en 1944.    en 1945.    en 1946.    en 1947.    en 1948.    en 1949.    en 1950.    en 1951.    en 1952.    en 1953.    en 1954.    en 1955.    en 1956.    en 1957.    en 1958.    en 1959.    en 1960.    en 1961.    en 1962.    en 1963.    en 1964.    en 1965.    en 1966.    en 1967.    en 1968.    en 1969.    en 1970.    en 1971.    en 1972.    en 1973.    en 1974.    en 1975.    en 1976.    en 1977.    en 1978.    en 1979.    en 1980.    en 1981.    en 1982.    en 1983.    en 1984.    en 1985.    en 1986.    en 1987.    en 1988.    en 1989.    en 1990.    en 1991.    en 1992.    en 1993.    en 1994.    en 1995.    en 1996.    en 1997.    en 1998.    en 1999.    en 2000.    en 2001.    en 2002.    en 2003.    en 2004.    en 2005.    en 2006.    en 2007.    en 2008.    en 2009.    en 2010.    en 2011.    en 2012.    en 2013.    en 2014.    en 2015.    en 2016.    en 2017.    en 2018.    en 2019.    en 2020.    en 2021.    en 2022.    en 2023.    en 2024.    en 2025.    en 2026.    en 2027.    en 2028.    en 2029.    en 2030.    en 2031.    en 2032.    en 2033.    en 2034.    en 2035.    en 2036.    en 2037.    en 2038.    en 2039.    en 2040.    en 2041.    en 2042.    en 2043.    en 2044.    en 2045.    en 2046.    en 2047.    en 2048.    en 2049.    en 2050.    en 2051.    en 2052.    en 2053.    en 2054.    en 2055.    en 2056.    en 2057.    en 2058.    en 2059.    en 2060.    en 2061.    en 2062.    en 2063.    en 2064.    en 2065.    en 2066.    en 2067.    en 2068.    en 2069.    en 2070.    en 2071.    en 2072.    en 2073.    en 2074.    en 2075.    en 2076.    en 2077.    en 2078.    en 2079.    en 2080.    en 2081.    en 2082.    en 2083.    en 2084.    en 2085.    en 2086.    en 2087.    en 2088.    en 2089.    en 2090.    en 2091.    en 2092.    en 2093.    en 2094.    en 2095.    en 2096.    en 2097.    en 2098.    en 2099.    en 2100.    en 2101.    en 2102.    en 2103.    en 2104.    en 2105.    en 2106.    en 2107.    en 2108.    en 2109.    en 2110.    en 2111.    en 2112.    en 2113.    en 2114.    en 2115.    en 2116.    en 2117.    en 2118.    en 2119.    en 2120.    en 2121.    en 2122.    en 2123.    en 2124.    en 2125.    en 2126.    en 2127.    en 2128.    en 2129.    en 2130.    en 2131.    en 2132.    en 2133.    en 2134.    en 2135.    en 2136.    en 2137.    en 2138.    en 2139.    en 2140.    en 2141.    en 2142.    en 2143.    en 2144.    en 2145.    en 2146.    en 2147.    en 2148.    en 2149.    en 2150.    en 2151.    en 2152.    en 2153.    en 2154.    en 2155.    en 2156.    en 2157.    en 2158.    en 2159.    en 2160.    en 2161.    en 2162.    en 2163.    en 2164.    en 2165.    en 2166.    en 2167.    en 2168.    en 2169.    en 2170.    en 2171.    en 2172.    en 2173.    en 2174.    en 2175.    en 2176.    en 2177.    en 2178.    en 2179.    en 2180.    en 2181.    en 2182.    en 2183.    en 2184.    en 2185.    en 2186.    en 2187.    en 2188.    en 2189.    en 2190.    en 2191.    en 2192.    en 2193.    en 2194.    en 2195.    en 2196.    en 2197.    en 2198.    en 2199.    en 2200.    en 2201.    en 2202.    en 2203.    en 2204.    en 2205.    en 2206.    en 2207.    en 2208.    en 2209.    en 2210.    en 2211.    en 2212.    en 2213.    en 2214.    en 2215.    en 2216.    en 2217.    en 2218.    en 2219.    en 2220.    en 2221.    en 2222.    en 2223.    en 2224.    en 2225.    en 2226.    en 2227.    en 2228.    en 2229.    en 2230.    en 2231.    en 2232.    en 2233.    en 2234.    en 2235.    en 2236.    en 2237.    en 2238.    en 2239.    en 2240.    en 2241.    en 2242.    en 2243.    en 2244.    en 2245.    en 2246.    en 2247.    en 2248.    en 2249.    en 2250.    en 2251.    en 2252.    en 2253.    en 2254.    en 2255.    en 2256.    en 2257.    en 2258.    en 2259.    en 2260.    en 2261.    en 2262.    en 2263.    en 2264.    en 2265.    en 2266.    en 2267.    en 2268.    en 2269.    en 2270.    en 2271.    en 2272.    en 2273.    en 2274.    en 2275.    en 2276.    en 2277.    en 2278.    en 2279.    en 2280.    en 2281.    en 2282.    en 2283.    en 2284.    en 2285.    en 2286.    en 2287.    en 2288.    en 2289.    en 2290.    en 2291.    en 2292.    en 2293.    en 2294.    en 2295.    en 2296.    en 2297.    en 2298.    en 2299.    en 2300.    en 2301.    en 2302.    en 2303.    en 2304.    en 2305.    en 2306.    en 2307.    en 2308.    en 2309.    en 2310.    en 2311.    en 2312.    en 2313.    en 2314.    en 2315.    en 2316.    en 2317.    en 2318.    en 2319.    en 2320.    en 2321.    en 2322.    en 2323.    en 2324.    en 2325.    en 2326.    en 2327.    en 2328.    en 2329.    en 2330.    en 2331.    en 2332.    en 2333.    en 2334.    en 2335.    en 2336.    en 2337.    en 2338.    en 2339.    en 2340.    en 2341.    en 2342.    en 2343.    en 2344.    en 2345.    en 2346.    en 2347.    en 2348.    en 2349.    en 2350.    en 2351.    en 2352.    en 2353.    en 2354.    en 2355.    en 2356.    en 2357.    en 2358.    en 2359.    en 2360.    en 2361.    en 2362.    en 2363.    en 2364.    en 2365.    en 2366.    en 2367.    en 2368.    en 2369.    en 2370.    en 2371.    en 2372.    en 2373.    en 2374.    en 2375.    en 2376.    en 2377.    en 2378.    en 2379.    en 2380.    en 2381.    en 2382.    en 2383.    en 2384.    en 2385.    en 2386.    en 2387.    en 2388.    en 2389.    en 2390.    en 2391.    en 2392.    en 2393.    en 2394.    en 2395.    en 2396.    en 2397.    en 2398.    en 2399.    en 2400.    en 2401.    en 2402.    en 2403.    en 2404.    en 2405.    en 2406.    en 2407.    en 2408.    en 2409.    en 2410.    en 2411.    en 2412.    en 2413.    en 2414.    en 2415.    en 2416.    en 2417.    en 2418.    en 2419.    en 2420.    en 2421.    en 2422.    en 2423.    en 2424.    en 2425.    en 2426.    en 2427.    en 2428.    en 2429.    en 2430.    en 2431.    en 2432.    en 2433.    en 2434.    en 2435.    en 2436.    en 2437.    en 2438.    en 2439.    en 2440.    en 2441.    en 2442.    en 2443.    en 2444.    en 2445.    en 2446.    en 2447.    en 2448.    en 2449.    en 2450.    en 2451.    en 2452.    en 2453.    en 2454.    en 2455.    en 2456.    en 2457.    en 2458.    en 2459.    en 2460.    en 2461.    en 2462.    en 2463.    en 2464.    en 2465.    en 2466.    en 2467.    en 2468.    en 2469.    en 2470.    en 2471.    en 2472.    en 2473.    en 2474.    en 2475.    en 2476.    en 2477.    en 2478.    en 2479.    en 2480.    en 2481.    en 2482.    en 2483.    en 2484.    en 2485.    en 2486.    en 2487.    en 2488.    en 2489.    en 2490.    en 2491.    en 2492.    en 2493.    en 2494.    en 2495.    en 2496.    en 2497.    en 2498.    en 2499.    en 2500.    en 2501.    en 2502.    en 2503.    en 2504.    en 2505.    en 2506.    en 2507.    en 2508.    en 2509.    en 2510.    en 2511.    en 2512.    en 2513.    en 2514.    en 2515.    en 2516.    en 2517.    en 2518.    en 2519.    en 2520.    en 2521.    en 2522.    en 2523.    en 2524.    en 2525.    en 2526.    en 2527.    en 2528.    en 2529.    en 2530.    en 2531.    en 2532.    en 2533.    en 2534.    en 2535.    en 2536.    en 2537.    en 2538.    en 2539.    en 2540.    en 2541.    en 2542.    en 2543.    en 2544.    en 2545.    en 2546.    en 2547.    en 2548.    en 2549.    en 2550.    en 2551.    en 2552.    en 2553.    en 2554.    en 2555.    en 2556.    en 2557.    en 2558.    en 2559.    en 2560.    en 2561.    en 2562.    en 2563.    en 2564.    en 2565.    en 2566.    en 2567.    en 2568.    en 2569.    en 2570.    en 2571.    en 2572.    en 2573.    en 2574.    en 2575.    en 2576.    en 2577.    en 2578.    en 2579.    en 2580.    en 2581.    en 2582.    en 2583.    en 2584.    en 2585.    en 2586.    en 2587.    en 2588.    en 2589.    en 2590.    en 2591.    en 2592.    en 2593.    en 2594.    en 2595.    en 2596.    en 2597.    en 2598.    en 2599.    en 2600.    en 2601.    en 2602.    en 2603.    en 2604.    en 2605.    en 2606.    en 2607.    en 2608.    en 2609.    en 2610.    en 2611.    en 2612.    en 2613.    en 2614.    en 2615.    en 2616.    en 2617.    en 2618.    en 2619.    en 2620.    en 2621.    en 2622.    en 2623.    en 2624.    en 2625.    en 2626.    en 2627.    en 2628.    en 2629.    en 2630.    en 2631.    en 2632.    en 2633.    en 2634.    en 2635.    en 2636.    en 2637.    en 2638.    en 2639.    en 2640.    en 2641.    en 2642.    en 2643.    en 2644.    en 2645.    en 2646.    en 2647.    en 2648.    en 2649.    en 2650.    en 2651.    en 2652.    en 2653.    en 2654.    en 2655.    en 2656.    en 2657.    en 2658.    en 2659.    en 2660.    en 2661.    en 2662.    en 2663.    en 2664.    en 2665.    en 2666.    en 2667.    en 2668.    en 2669.    en 2670.    en 2671.    en 2672.    en 2673.    en 2674.    en 2675.    en 2676.    en 2677.    en 2678.    en 2679.    en 2680.    en 2681.    en 2682.    en 2683.    en 2684.    en 2685.    en 2686.    en 2687.    en 2688.    en 2689.    en 2690.    en 2691.    en 2692.    en 2693.    en 2694.    en 2695.    en 2696.    en 2697.    en 2698.    en 2699.    en 2700.    en 2701.    en 2702.    en 2703.    en 2704.    en 2705.    en 2706.    en 2707.    en 2708.    en 2709.    en 2710.    en 2711.    en 2712.    en 2713.    en 2714.    en 2715.    en 2716.    en 2717.    en 2718.    en 2719.    en 2720.    en 2721.    en 2722.    en 2723.    en 2724.    en 2725.    en 2726.    en 2727.    en 2728.    en 2729.    en 2730.    en 2731.    en 2732.    en 2733.    en 2734.    en 2735.    en 2736.    en 2737.    en 2738.    en 2739.    en 2740.    en 2741.    en 2742.    en 2743.    en 2744.    en 2745.    en 2746.    en 2747.    en 2748.    en 2749.    en 2750.    en 2751.    en 2752.    en 2753.    en 2754.    en 2755.    en 2756.    en 2757.    en 2758.    en 2759.    en 2760.    en 2761.    en 2762.    en 2763.    en 2764.    en 2765.    en 2766.    en 2767.    en 2768.    en 2769.    en 2770.    en 2771.    en 2772.    en 2773.    en 2774.    en 2775.    en 2776.    en 2777.    en 2778.    en 2779.    en 2780.    en 2781.    en 2782.    en 2783.    en 2784.    en 2785.    en 2786.    en 2787.    en 2788.    en 2789.    en 2790.    en 2791.    en 2792.    en 2793.    en 2794.    en 2795.    en 2796.    en 2797.    en 2798.    en 2799.    en 2800.    en 2801.    en 2802.    en 2803.    en 2804.    en 2805.    en 2806.    en 2807.    en 2808.    en 2809.    en 2810.    en 2811.    en 2812.    en 2813.    en 2814.    en 2815.    en 2816.    en 2817.    en 2818.    en 2819.    en 2820.    en 2821.    en 2822.    en 2823.    en 2824.    en 2825.    en 2826.    en 2827.    en 2828.    en 2829.    en 2830.    en 2831.    en 2832.    en 2833.    en 2834.    en 2835.    en 2836.    en 2837.    en 2838.    en 2839.    en 2840.    en 2841.    en 2842.    en 2843.    en 2844.    en 2845.    en 2846.    en 2847.    en 2848.    en 2849.    en 2850.    en 2851.    en 2852.    en 2853.    en 2854.    en 2855.    en 2856.    en 2857.    en 2858.    en 2859.    en 2860.    en 2861.    en 2862.    en 2863.    en 2864.    en 2865.    en 2866.    en 2867.    en 2868.    en 2869.    en 2870.    en 2871.    en 2872.    en 2873.    en 2874.    en 2875.    en 2876.    en 2877.    en 2878.    en 2879.    en 2880.    en 2881.    en 2882.    en 2883.    en 2884.    en 2885.    en 2886.    en 2887.    en 2888.    en 2889.    en 2890.    en 2891.    en 2892.    en 2893.    en 2894.    en 2895.    en 2896.    en 2897.    en 2898.    en 2899.    en 2900.    en 2901.    en 2902.    en 2903.    en 2904.    en 2905.    en 2906.    en 2907.    en 2908.    en 2909.    en 2910.    en 2911.    en 2912.    en 2913.    en 2914.    en 2915.    en 2916.    en 2917.    en 2918.    en 2919.    en 2920.    en 2921.    en 2922.    en 2923.    en 2924.    en 2925.    en 2926.    en 2927.    en 2928.    en 2929.    en 2930.    en 2931.    en 2932.    en 2933.    en 2934.    en 2935.    en 2936.    en 2937.    en 2938.    en 2939.    en 2940.    en 2941.    en 2942.    en 2943.    en 2944.    en 2945.    en 2946.    en 2947.    en 2948.    en 2949.    en 2950.    en 2951.    en 2952.    en 2953.    en 2954.    en 2955.    en 2956.    en 2957.    en 2958.    en 2959.    en 2960.    en 2961.    en 2962.    en 2963.    en 2964.    en 2965.    en 2966.    en 2967.    en 2968.    en 2969.    en 2970.    en 2971.    en 2972.    en 2973.    en 2974.    en 2975.    en 2976.    en 2977.    en 2978.    en 2979.    en 2980.    en 2981.    en 2982.    en 2983.    en 2984.    en 2985.    en 2986.    en 2987.    en 2988.    en 2989.    en 299

konden ze ook worden omgezet in FoLiA met het programma dat was gebruikt voor de recente tekst<sup>4</sup>.

## 4 Verwerking met Frog

Nadat de drie teksten waren omgezet in het formaat FoLiA, konden ze zonder probleem binnen enkele minuten worden verwerkt door Frog. Het programma genereerde drie nieuwe FoLiA-bestanden met daarin de taalkundige analyse van de tekst in elk document. Voor de vergelijking gebruiken we onze eigen FoLiA-browser, een verzameling van Javascriptprogramma's voor visualisatie van de inhoud van FoLiA-bestanden.

## 5 Scankwaliteit

We begonnen onze analyse met de controle van de kwaliteit van de teksten. De twee oude teksten zijn gedigitaliseerd met behulp van optical character recognition (ocr) en hierdoor zijn sommige letters verkeerd herkend. Van de tweede bladzijde van elk document telden we de woorden en getallen (niet de leestekens), en controleerden we met behulp van de PDF-bestanden op de websites hoeveel daarvan correct waren herkend:

- Hardy (1750): 139 woorden; 2 fouten, 98% correct
- Stuart (1800): 147 woorden; 4 fouten, 97% correct
- Wiesbaden (1993): 289 woorden: 0 fouten, 100% correct

Bij deze controle hebben we herkenning van de lange s (f) als de letter f goedgerekend. Deze verwisseling kwam diverse keren voor in de twee oude documenten en zou met behulp van een woordenlijst gemakkelijk te corrigeren moeten zijn. Voor de gecontroleerde bladzijden vallen de foutaantallen erg mee.

---

<sup>4</sup>Ko van der Sloot van de Universiteit Tilburg heeft ook een programma geschreven voor het omzetten van EDBO-documenten naar FoLiA-formaat. Het extra materiaal dat in deze sectie wordt gemeld komt bij die documenten ook in de FoLiA-versie terecht.

## 6 Zinsgrenzen

Zinsgrenzen zijn in de oudere documenten soms lastig te herkennen omdat niet consequent gebruik wordt gemaakt van leestekens:

- Hardy (1750): 9 zinnen; 9 fouten, 0% correct
- Stuart (1800): 4 zinnen; 0 fouten, 100% correct
- Wiesbaden (1993): 14 zinnen: 0 fouten, 100% correct

Het programma heeft alleen bij de tekst van Hary moeite om de zinsgrenzen te vinden. Dit komt doordat deze tekst in dichtvorm is opgezet met gebruik van komma's en inspringen om zinsgrenzen aan te geven. Frog herkent hier dan ook maar een zinsgrens, op een verkeerde plek door een vraagteken die is ontstaan door een ocr-fout.

## 7 Lemma's

Frog bepaalt voor elk woord in een tekst het lemma, wat kan worden gezien als de corresponderende woordenboekingang. Voor veel woordvormen is maar een mogelijke woordenboekingang. Voor andere woorden moet een keuze gemaakt worden. Zo zou de woordvorm *vis* zowel kunnen passen bij het gelijkvormige zelfstandige naamwoord als bij het werkwoord *vissen*.

- Hardy (1750): 139 woorden; 38 fouten, 72% correct
- Stuart (1800): 147 woorden; 23 fouten, 84% correct
- Wiesbaden (1993): 289 woorden: 6 fouten, 98% correct

In het evaluatieproces hebben we ocrfouten op basis van de verwisselingen lange  $s \leftrightarrow f$  en  $ij \leftrightarrow y$  geaccepteerd. Verder hebben we van alle woorden vereist dat zij werden gelinkt aan een correct hedendaags equivalent. Het lemmatisatieproces genereert meer fouten bij het verwerken van oudere teksten. De fouten worden voornamelijk veroorzaakt door vocabulaireverschillen en ocrfouten.

## 8 Syntactische woordklassen (part-of-speech)

Frog kent ook syntactische klassen als werkwoord en zelfstandig naamwoord toe aan woorden. Bij de evaluatie hebben we alleen naar de hoofdklassen gekeken. Eigenschappen als enkelvoud vs meervoud en tijdsbepalingen hebben we buiten beschouwing gelaten.

- Hardy (1750): 139 woorden; 38 fouten, **72% correct**
- Stuart (1800): 147 woorden; 24 fouten, **84% correct**
- Wiesbaden (1993): 289 woorden: 16 fouten, **94% correct**

De scores zijn vergelijkbaar met die van de lemma-analyse. Part-of-speech tags worden in het hedendaagse werk slechter herkend dan lemma's. De meest voorkomende fout is het benoemen van bijwoorden als bijvoeglijke naamwoorden.

## 9 Namen (named entities)

Frog bevat ook een module voor naamherkenning. Naast het onderscheid tussen namen en andere woorden, wordt ook de klasse van een herkende naam aangegeven, bijvoorbeeld: persoon, organisatie of locatie. We hebben alleen het onderscheid wel-geen getest:

- Hardy (1750): 4 namen; 2 correct, 14 fouten, **13% correct**
- Stuart (1800): 8 woorden; 3 correct, 11 fouten, **21% correct**
- Wiesbaden (1993): 4 namen: 4 correct, 0 fouten, **100% correct**

De namen in de hedendaagse tekst worden correct herkend. Voor de oudere teksten maakt het systeem meer fouten dan dat het namen correct herkent. Een probleem hierbij is capitalisatie: in de oude teksten is de eerste letter van diverse zelfstandige woorden een hoofdletter. Daarnaast bevatte een tekst (Stuart) namen in smallcaps, wat door de optical character recognition werd omgezet kleine letters en daarna lastig was te herkennen als naam.

## 10 Conclusies

We hebben twee oude teksten verwerkt met het taalanalyseprogramma Frog en hebben vervolgens de resultaten vergeleken met die van een recente tekst. Zoals verwacht, werden de oude teksten slechter verwerkt dan de recente tekst. De prestaties van het systeem waren het slechtst voor de oudste tekst en voor complexere analyses. Het herkennen van zinsgrenzen ging, afhankelijk van de aangeboden tekst, goed. Herkenning van lemma's, syntactische woordklassen en namen ging beduidend slechter bij de twee oude teksten.

Bij de resultaten moet worden aangetekend worden dat we voor de test niet de meest ingewikkelde bladzijden hebben uitgekozen. Daarnaast hebben we extra materiaal, zoals paginanummers, kopjes en voetnoten, handmatig van de bladzijden verwijderd. Als we deze informatie in de bestanden hadden laten staan dan waren de testresultaten ongetwijfeld slechter geweest. Het opschonen van de teksten is wenselijk maar het is ondoenlijk om dit handmatig te doen voor al ons materiaal en het nog maar de vraag of dit proces automatiseerbaar is.

We kunnen op dit moment de EDBO-documenten verwerken met Frog maar de kwaliteit van de resultaten zal niet erg hoog zijn. Voor een verbetering van de analyseresultaten kunnen de volgende processen proberen te verbeteren:

1. **Tekenherkenning:** verbeterde letterherkenning (optical character recognition) zal waarschijnlijk leiden tot een verbetering van alle analyses van Frog. Het ocr-proces kunnen we niet overdoen maar mogelijk kan postprocessing (bijvoorbeeld met TICCL of een lijst van veelgemaakte fouten) de kwaliteit van de teksten verbeteren.
2. **Markering van buitentekstelijk materiaal:** paginanummers, kopjes en noten moeten apart worden verwerkt door het programma. Als we deze al gemarkeerd zijn of als zij automatisch kunnen worden herkend dan kunnen daardoor de analyses van Frog worden verbeterd.
3. **Toevoeging lexicon ouder Nederlands:** dit zou helpen om van meer woorden het lemma en de woordklasse goed te herkennen. Het INL kan mogelijk geschikte lexica leveren.
4. **Recapitalisatie:** omdat de herkenning van namen sterkt leunt op de aanwezigheid en afwezigheid van hoofdletters, heeft het alternatieve hoofdlettergebruik in oude teksten een negatieve invloed op het herkennen van namen. Vooraf het hoofdlettergebruik standaardiseren zou de naamherkenning verbeteren.

Niet voor alle taken is kant-en-klare software beschikbaar. Verbetering van de resultaten van Frog op oude Nederlandse teksten zal extra werk kosten maar hier is bij de planning van Nederlab rekening mee gehouden.

## A Geteste teksten (zoals gescand)

### Hardy (1750)

Uw vlugge ftyl , hoewel \* vvat kreupel afgefchreyen , Gaf ftraks een ftaaltje van den hersfen-ryken bol , Ik zag geleertheid en taalkunde daar in leven , Van Godtvrugt en verftandt en leesvrugt even vol : My dagt , ik zag U naar den groten Tempel flappen , Vol geest en vuur , verzelt van ? t agtbaar Priesterdom \* t Geen U geleidde naar dehooge kansfel trappen , Den aandagt ftelen van den Godtgewyden drom : My dagt , ik hoorde daar de blyde Maagden reij en In 't feestkleed uitgedost , voor uwe voeten neer Gezeten , vol van vreugd haar heilig danklied fprijen , Ten prys van U , maar ook vooral van Uwen Heer , 'k Meen ' Jezus , die Uw hart door Zynen gloed deet branden , Uw tong ontftekende met heilig autervuur Niet meer van ftieren vleesch ; maar beetere offerhanden , In vollen vlam gezet , niet binnen Zalems muur , Maar in de vrugtbare en genaderyke ftreken Van Neetlands Zion , daar de vrede Koning woont ,

### Stuart (1800)

vooral den Romeinfchen Staat . Naauwlijks had het openbaar geweld van het op nieuw verbondene Driemanfchap aan Rome in crassus en pompejus Bewindsluiden opgedrongen , die Hechts den fchijrt van een vrij Gemeenebest in wezen zouden laten , of cato , een hardnekkige verdediger der burgerlijke vrijheid , hervattede den wanhoopigen ftrijd , zonder aan de onverfchilligheid , infchiklijkheid of zwakheid van anderen eene te duur betaalde rust te vergunnen . , Even min door de behaalde zege zijner tegenpartije , als door zijn naauwlijks ontkomen lijfsgevaar ( 1 ) afgefchrikt van eene " nieuwe en ftoute pooging , tradt cato ) moedig als mededinger voor naar het Pree [ torfchap van dit jaar , met geen minder oogmerk , dan om , van agter dit gezag verfchanst , de vijanden van het Gemeenebest meer op eene gelijke hoogte te beftrijden ( 2 ) . De Confuls , wier werk het zijn moest , de verkiezing der overige Overheden voor het reeds aangevangen re -

## Wiesbaden (1993)

X-en had weer zo'n origineel thema ! Samengevat zou men kunnen stellen dat ze a.h.w. scrabbelen met hetzelfde alfabet en dezelfde woord - en letterwaarden ....

SPELEN betekent oorspronkelijk 'zich continu bewegen , zich vlot kunnen bewegen . Het spel wordt daartoe beschermd door spelregels . Speelt men , dan moet men dus ernstig spelen . Overtredingen gelden enkel binnen het spel . Maar wie weigert het spel mee te spelen , bekritiseert niet zozeer de manier waarop het spel gespeeld wordt , maar het spel zelf , en de overtuiging , het geloof en de wil dat er de basis van vormt . En dat is voor de spelers die het paradoxaliter ernstig menen , onvergeeflijk , en vervult de sfeer met ongemak , wrevel tot agressie . Om een 'verbroken spel ' te vermijden , geldt dus de eis van de algemene instemming , zeg maar ' conformisme ' .

Als men dus zo opgaat in het spel , dat het werkelijkheid wordt , dan ' speelt ' men niet meer in de eigenlijke betekenis van het woord . Misschien bestaat de inwijding in l'art de vivre dan ook in het vlotten van de evenwichtsoefeningen tussen spelen en leven , die zich met al onze menselijke activiteiten vermengd hebben : het spel van de lektuur , van het gasten-ontvangen , van het reizen , van het zakendoen , ...

Eenzijds ligt het belang van het spel(en) in de mogelijkheid te ontsnappen aan de vernauwde wereld van het zakelijke , van de orde van de noodzakelijkheden , de mens te bevrijden van het determinisme , om zo meer en meer zichzelf te vinden .

Anderzijds , zoals verder nog ter sprake zal komen , ontstaan via het spelen nieuwe determinanten , wordt de ' onvrijheid ' in de hand gewerkt , en verliest de mens uiteindelijk juist zijn kans op persoonlijkheid .

Zo kan men zich inbeelden dat Benoît's personages sculpturen als burens hebben , Van Beirendoncks creaties appreciëren maar liever Armani dragen , Wittamer de ontdekking

## Referenties

- [1] A. Van den Bosch, G.J. Busser, W. Daelemans, and S. Canisius. An efficient memory-based morphosyntactic tagger and parser for Dutch. In *Selected Papers of the 17th Computational Linguistics in the Netherlands Meeting*, pages 99–114. Leuven, Belgium, 2007.



- [2] M. van Gompel and M. Reynaert. FoLiA: A practical XML format for linguistic annotation - a descriptive and comparative study. *CLIN Journal*, 3:63–81, 2013.