

|                              |      | Average improvement in F1 scores |       |       |        |       |       |       |       |
|------------------------------|------|----------------------------------|-------|-------|--------|-------|-------|-------|-------|
| Percentage of improved cases | n=1  | Y                                | YP    | YT    | YW     | YPT   | YPW   | YTW   | YPTW  |
|                              | Y    |                                  | 0.3%  | 2.2%  | 0.3%   | -0.4% | 0.3%  | 2.2%  | -0.4% |
|                              | YP   | 58.8%                            |       | 9.9%  | -0.3%  | 7.0%  | 0.0%  | 9.9%  | 7.0%  |
|                              | YT   | 100.0%                           | 90.9% |       | -2.1%  | -2.5% | -9.0% | 0.0%  | -2.5% |
|                              | YW   | 0.0%                             | 41.2% | 0.0%  |        | -0.4% | 0.3%  | 2.2%  | -0.4% |
|                              | YPT  | 72.7%                            | 90.9% | 54.5% | 72.7%  |       | -6.6% | 2.6%  | 0.0%  |
|                              | YPW  | 58.8%                            | 0.0%  | 9.1%  | 58.8%  | 9.1%  |       | 9.9%  | 7.0%  |
|                              | YTW  | 100.0%                           | 90.9% | 0.0%  | 100.0% | 45.5% | 90.9% |       | -2.5% |
|                              | YPTW | 72.7%                            | 90.9% | 54.5% | 72.7%  | 0.0%  | 90.9% | 54.5% |       |