# Symptom

How to reproduce?

➢ Interrupt-heavy workloads: YCSB, iPerf, etc.
➢ Bind IRQs to a specific socket/NUMA node
    ➢ Network performance is very sensitive to having IRQs routed to the "wrong" socket because a PCI bus is usually connected to one socket. Some even reported up to [2x](#) slower performance

Tasks are constantly getting pulled to the socket/NUMA node that IRQs are bound to while leaving other sockets nearly idle.

➢ Within each socket, loads are fairly balanced
➢ Spreading out tasks more **evenly** across sockets can improve performance numbers from YCSB benchmark under light load.

# Light Load

```
06:18:58 PM  CPU   %usr  %nice  %sys  %iowait  %irq  %soft  %steal  %guest  %gnice  %idle
06:19:01 PM  all  18.00  0.00   5.49   0.00   0.00  1.91   0.00    0.00    0.00   74.60
06:19:01 PM    0  21.11  0.00   7.41   0.00   0.00  0.74   0.00    0.00    0.00   70.74
06:19:01 PM    1  21.00  0.00   5.69   0.00   0.00  0.71   0.00    0.00    0.00   72.60
06:19:01 PM    2  20.64  0.00   6.76   0.00   0.00  1.42   0.00    0.00    0.00   71.17
06:19:01 PM    3  20.64  0.00   6.05   0.00   0.00  0.71   0.00    0.00    0.00   72.60
06:19:01 PM    4  19.06  0.00   6.47   0.00   0.00  1.80   0.00    0.00    0.00   72.66
06:19:01 PM    5  21.22  0.00   5.76   0.00   0.00  1.44   0.00    0.00    0.00   71.58
06:19:01 PM    6  21.38  0.00   6.88   0.00   0.00  2.54   0.00    0.00    0.00   69.20
06:19:01 PM    7  19.08  0.00   5.30   0.00   0.00  0.00   0.00    0.00    0.00   75.62
06:19:01 PM    8  19.38  0.00   5.88   0.00   0.00  0.00   0.00    0.00    0.00   74.74
06:19:01 PM    9  19.30  0.00   5.61   0.00   0.00  0.00   0.00    0.00    0.00   75.09
06:19:01 PM   10  18.69  0.00   6.57   0.00   0.00  0.00   0.00    0.00    0.00   74.74
06:19:01 PM   11  19.86  0.00   5.23   0.00   0.00  0.00   0.00    0.00    0.00   74.91
06:19:01 PM   12  18.75  0.00   5.90   0.00   0.00  0.00   0.00    0.00    0.00   75.35
06:19:01 PM   13  19.10  0.00   5.21   0.00   0.00  0.00   0.00    0.00    0.00   75.69
06:19:01 PM   14  17.77  0.00   6.27   0.00   0.00  0.00   0.00    0.00    0.00   75.96
06:19:01 PM   15  47.45  0.00  10.98   0.00   0.00  23.14  0.00    0.00    0.00   18.43
06:19:01 PM   16  48.89  0.00  14.07   0.00   0.00  18.89  0.00    0.00    0.00   18.15
06:19:01 PM   17  50.55  0.00  13.65   0.00   0.00  21.40  0.00    0.00    0.00   14.39
06:19:01 PM   18  49.82  0.00  14.08   0.00   0.00  24.91  0.00    0.00    0.00   11.19
06:19:01 PM   19  53.45  0.00  14.18   0.00   0.00  20.00  0.00    0.00    0.00   12.36
06:19:01 PM   20  52.35  0.00  12.64   0.00   0.00  25.99  0.00    0.00    0.00    9.03
06:19:01 PM   21  52.71  0.00  13.72   0.00   0.00  23.83  0.00    0.00    0.00    9.75
06:19:01 PM   22  52.55  0.00  13.50   0.00   0.00  25.18  0.00    0.00    0.00    8.76
06:19:01 PM   23  56.54  0.00  18.02   0.00   0.00  1.41   0.00    0.00    0.00   24.03
06:19:01 PM   24   1.68  0.00   1.01   0.00   0.00  0.00   0.00    0.00    0.00   97.32
06:19:01 PM   25   0.67  0.00   0.67   0.00   0.00  0.00   0.00    0.00    0.00   98.66
06:19:01 PM   26   0.34  0.00   0.00   0.00   0.00  0.00   0.00    0.00    0.00   99.66
06:19:01 PM   27   0.00  0.00   0.00   0.00   0.00  0.00   0.00    0.00    0.00  100.00
06:19:01 PM   28   1.67  0.00   0.33   0.00   0.00  0.00   0.00    0.00    0.00   97.99
06:19:01 PM   29   0.67  0.00   0.33   0.00   0.00  0.00   0.00    0.00    0.00   99.00
06:19:01 PM   30   0.34  0.00   0.00   0.00   0.00  0.00   0.00    0.00    0.00   99.66
06:19:01 PM   31   1.00  0.00   0.33   0.00   0.00  0.00   0.00    0.00    0.00   98.66
06:19:01 PM   32   1.00  0.00   0.67   0.00   0.00  0.00   0.00    0.00    0.00   98.33
06:19:01 PM   33   1.00  0.00   0.67   0.00   0.00  0.00   0.00    0.00    0.00   98.33
06:19:01 PM   34   1.34  0.00   1.00   0.00   0.00  0.00   0.00    0.00    0.00   97.66
06:19:01 PM   35   2.34  0.00   4.35   0.00   0.00  0.00   0.00    0.00    0.00   93.31
06:19:01 PM   36   3.37  0.00   2.69   0.00   0.00  0.00   0.00    0.00    0.00   93.94
06:19:01 PM   37   1.33  0.00   0.33   0.00   0.00  0.00   0.00    0.00    0.00   98.33
06:19:01 PM   38   1.34  0.00   0.33   0.00   0.00  0.00   0.00    0.00    0.00   98.33
06:19:01 PM   39   1.00  0.00   0.67   0.00   0.00  0.00   0.00    0.00    0.00   98.33
06:19:01 PM   40   0.67  0.00   0.33   0.00   0.00  0.00   0.00    0.00    0.00   99.00
06:19:01 PM   41   0.66  0.00   1.00   0.00   0.00  0.00   0.00    0.00    0.00   98.34
06:19:01 PM   42   1.00  0.00   0.67   0.00   0.00  0.00   0.00    0.00    0.00   98.33
06:19:01 PM   43   0.67  0.00   0.33   0.00   0.00  0.00   0.00    0.00    0.00   99.00
06:19:01 PM   44   1.34  0.00   1.68   0.00   0.00  0.00   0.00    0.00    0.00   96.98
06:19:01 PM   45   0.00  0.00   0.33   0.00   0.00  0.00   0.00    0.00    0.00   99.67
06:19:01 PM   46   0.33  0.00   0.33   0.00   0.00  0.00   0.00    0.00    0.00   99.34
06:19:01 PM   47   0.00  0.00   0.33   0.00   0.00  0.00   0.00    0.00    0.00   99.67
```

```
06:19:01 PM   48  55.32  0.00  18.44   0.00   0.00  0.00   0.00    0.00    0.00   26.24
06:19:01 PM   49  56.49  0.00  17.19   0.00   0.00  0.00   0.00    0.00    0.00   26.32
06:19:01 PM   50  55.63  0.00  16.55   0.00   0.00  0.00   0.00    0.00    0.00   27.82
06:19:01 PM   51  54.74  0.00  16.14   0.00   0.00  0.00   0.00    0.00    0.00   29.12
06:19:01 PM   52  53.15  0.00  17.13   0.00   0.00  0.00   0.00    0.00    0.00   29.72
06:19:01 PM   53  53.19  0.00  14.54   0.00   0.00  0.00   0.00    0.00    0.00   32.27
06:19:01 PM   54  50.90  0.00  14.80   0.00   0.00  0.00   0.00    0.00    0.00   34.30
06:19:01 PM   55  49.64  0.00  14.86   0.00   0.00  0.00   0.00    0.00    0.00   35.51
06:19:01 PM   56  47.62  0.00  14.29   0.00   0.00  0.00   0.00    0.00    0.00   38.10
06:19:01 PM   57  46.93  0.00  14.08   0.00   0.00  0.00   0.00    0.00    0.00   38.99
06:19:01 PM   58  45.36  0.00  13.57   0.00   0.00  0.00   0.00    0.00    0.00   41.07
06:19:01 PM   59  43.93  0.00  13.21   0.00   0.00  0.00   0.00    0.00    0.00   42.86
06:19:01 PM   60  42.20  0.00  12.77   0.00   0.00  0.00   0.00    0.00    0.00   45.04
06:19:01 PM   61  40.66  0.00  11.72   0.00   0.00  0.00   0.00    0.00    0.00   47.62
06:19:01 PM   62  38.85  0.00  11.87   0.00   0.00  0.00   0.00    0.00    0.00   49.28
06:19:01 PM   63  35.46  0.00   9.57   0.00   0.00  0.00   0.00    0.00    0.00   54.96
06:19:01 PM   64  33.22  0.00   9.89   0.00   0.00  0.00   0.00    0.00    0.00   56.89
06:19:01 PM   65  31.10  0.00   9.19   0.00   0.00  0.00   0.00    0.00    0.00   59.72
06:19:01 PM   66  28.93  0.00   7.86   0.00   0.00  0.00   0.00    0.00    0.00   63.21
06:19:01 PM   67  28.83  0.00   8.54   0.00   0.00  0.00   0.00    0.00    0.00   62.63
06:19:01 PM   68  26.69  0.00   8.19   0.00   0.00  0.00   0.00    0.00    0.00   65.12
06:19:01 PM   69  25.87  0.00   7.69   0.00   0.00  0.00   0.00    0.00    0.00   66.43
06:19:01 PM   70  23.43  0.00   6.64   0.00   0.00  0.00   0.00    0.00    0.00   69.93
06:19:01 PM   71  19.08  0.00   7.07   0.00   0.00  0.00   0.00    0.00    0.00   73.85
06:19:01 PM   72   0.33  0.00   0.33   0.00   0.00  0.00   0.00    0.00    0.00   99.33
06:19:01 PM   73   0.67  0.00   0.33   0.00   0.00  0.00   0.00    0.00    0.00   99.00
06:19:01 PM   74   0.33  0.00   0.33   0.00   0.00  0.00   0.00    0.00    0.00   99.33
06:19:01 PM   75   1.67  0.00   0.33   0.00   0.00  0.00   0.00    0.00    0.00   98.00
06:19:01 PM   76   0.34  0.00   0.00   0.00   0.00  0.00   0.00    0.00    0.00   99.66
06:19:01 PM   77   0.33  0.00   0.33   0.00   0.00  0.00   0.00    0.00    0.00   99.33
06:19:01 PM   78   0.33  0.00   0.33   0.00   0.00  0.00   0.00    0.00    0.00   99.33
06:19:01 PM   79   1.33  0.00   1.00   0.00   0.00  0.00   0.00    0.00    0.00   97.67
06:19:01 PM   80   0.67  0.00   0.34   0.00   0.00  0.00   0.00    0.00    0.00   98.99
06:19:01 PM   81   0.67  0.00   0.00   0.00   0.00  0.00   0.00    0.00    0.00   99.33
06:19:01 PM   82   1.00  0.00   0.67   0.00   0.00  0.00   0.00    0.00    0.00   98.33
06:19:01 PM   83   0.67  0.00   0.33   0.00   0.00  0.00   0.00    0.00    0.00   99.00
06:19:01 PM   84   0.34  0.00   0.34   0.00   0.00  0.00   0.00    0.00    0.00   99.32
06:19:01 PM   85   0.33  0.00   0.33   0.00   0.00  0.00   0.00    0.00    0.00   99.33
06:19:01 PM   86   0.33  0.00   0.00   0.00   0.00  0.00   0.00    0.00    0.00   99.67
06:19:01 PM   87   1.67  0.00   1.67   0.00   0.00  0.00   0.00    0.00    0.00   96.67
06:19:01 PM   88   1.00  0.00   0.67   0.00   0.00  0.00   0.00    0.00    0.00   98.33
06:19:01 PM   89   1.00  0.00   0.67   0.00   0.00  0.00   0.00    0.00    0.00   99.33
06:19:01 PM   90   0.33  0.00   0.33   0.00   0.00  0.00   0.00    0.00    0.00   99.33
06:19:01 PM   91   0.00  0.00   0.00   0.00   0.00  0.00   0.00    0.00    0.00  100.00
06:19:01 PM   92   0.67  0.00   0.00   0.00   0.00  0.00   0.00    0.00    0.00   99.33
06:19:01 PM   93   2.66  0.00   2.66   0.00   0.00  0.00   0.00    0.00    0.00   94.68
06:19:01 PM   94   0.67  0.00   0.67   0.00   0.00  0.00   0.00    0.00    0.00   98.67
06:19:01 PM   95   0.33  0.00   0.33   0.00   0.00  0.00   0.00    0.00    0.00   99.33
```

# Heavy Load

| time | CPU | %usr | %nice | %sys | %iowait | %irq | %soft | %steal | %guest | %gnice | %idle |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 06:24:37 PM | all | 49.03 | 0.00 | 8.46 | 0.00 | 0.00 | 3.34 | 0.00 | 0.00 | 0.00 | 0.00 | 39.17 |
| 06:24:37 PM | 0 | 67.93 | 0.00 | 11.03 | 0.00 | 0.00 | 4.83 | 0.00 | 0.00 | 0.00 | 0.00 | 16.21 |
| 06:24:37 PM | 1 | 69.34 | 0.00 | 9.76 | 0.00 | 0.00 | 4.53 | 0.00 | 0.00 | 0.00 | 0.00 | 16.38 |
| 06:24:37 PM | 2 | 68.06 | 0.00 | 11.11 | 0.00 | 0.00 | 3.82 | 0.00 | 0.00 | 0.00 | 0.00 | 17.01 |
| 06:24:37 PM | 3 | 69.73 | 0.00 | 9.52 | 0.00 | 0.00 | 4.08 | 0.00 | 0.00 | 0.00 | 0.00 | 16.67 |
| 06:24:37 PM | 4 | 68.62 | 0.00 | 10.69 | 0.00 | 0.00 | 4.14 | 0.00 | 0.00 | 0.00 | 0.00 | 16.55 |
| 06:24:37 PM | 5 | 69.31 | 0.00 | 10.00 | 0.00 | 0.00 | 3.79 | 0.00 | 0.00 | 0.00 | 0.00 | 16.90 |
| 06:24:37 PM | 6 | 67.13 | 0.00 | 11.42 | 0.00 | 0.00 | 4.15 | 0.00 | 0.00 | 0.00 | 0.00 | 17.30 |
| 06:24:37 PM | 7 | 69.55 | 0.00 | 12.11 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 18.34 |
| 06:24:37 PM | 8 | 71.09 | 0.00 | 10.88 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 18.03 |
| 06:24:37 PM | 9 | 69.73 | 0.00 | 11.90 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 18.37 |
| 06:24:37 PM | 10 | 70.65 | 0.00 | 11.60 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 17.75 |
| 06:24:37 PM | 11 | 70.55 | 0.00 | 11.30 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 18.15 |
| 06:24:37 PM | 12 | 70.10 | 0.00 | 11.68 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 18.21 |
| 06:24:37 PM | 13 | 70.55 | 0.00 | 11.30 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 18.15 |
| 06:24:37 PM | 14 | 69.97 | 0.00 | 11.60 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 18.43 |
| 06:24:37 PM | 15 | 53.56 | 0.00 | 9.83 | 0.00 | 0.00 | 34.92 | 0.00 | 0.00 | 0.00 | 0.00 | 1.69 |
| 06:24:37 PM | 16 | 55.74 | 0.00 | 8.45 | 0.00 | 0.00 | 33.78 | 0.00 | 0.00 | 0.00 | 0.00 | 2.03 |
| 06:24:37 PM | 17 | 53.87 | 0.00 | 7.74 | 0.00 | 0.00 | 37.04 | 0.00 | 0.00 | 0.00 | 0.00 | 1.35 |
| 06:24:37 PM | 18 | 53.22 | 0.00 | 7.46 | 0.00 | 0.00 | 37.97 | 0.00 | 0.00 | 0.00 | 0.00 | 1.36 |
| 06:24:37 PM | 19 | 55.41 | 0.00 | 8.11 | 0.00 | 0.00 | 35.14 | 0.00 | 0.00 | 0.00 | 0.00 | 1.35 |
| 06:24:37 PM | 20 | 56.27 | 0.00 | 7.80 | 0.00 | 0.00 | 34.58 | 0.00 | 0.00 | 0.00 | 0.00 | 1.36 |
| 06:24:37 PM | 21 | 56.23 | 0.00 | 8.42 | 0.00 | 0.00 | 34.01 | 0.00 | 0.00 | 0.00 | 0.00 | 1.35 |
| 06:24:37 PM | 22 | 55.70 | 0.00 | 7.72 | 0.00 | 0.00 | 35.23 | 0.00 | 0.00 | 0.00 | 0.00 | 1.34 |
| 06:24:37 PM | 23 | 73.04 | 0.00 | 12.29 | 0.00 | 0.00 | 3.75 | 0.00 | 0.00 | 0.00 | 0.00 | 10.92 |
| 06:24:37 PM | 24 | 30.93 | 0.00 | 5.84 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 63.23 |
| 06:24:37 PM | 25 | 28.52 | 0.00 | 6.19 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 65.29 |
| 06:24:37 PM | 26 | 29.69 | 0.00 | 5.46 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 64.85 |
| 06:24:37 PM | 27 | 28.42 | 0.00 | 6.51 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 65.07 |
| 06:24:37 PM | 28 | 29.31 | 0.00 | 5.52 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 65.17 |
| 06:24:37 PM | 29 | 29.49 | 0.00 | 6.44 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 64.07 |
| 06:24:37 PM | 30 | 29.21 | 0.00 | 5.15 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 65.64 |
| 06:24:37 PM | 31 | 29.55 | 0.00 | 5.15 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 65.29 |
| 06:24:37 PM | 32 | 27.99 | 0.00 | 6.83 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 65.19 |
| 06:24:37 PM | 33 | 29.21 | 0.00 | 5.84 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 64.95 |
| 06:24:37 PM | 34 | 28.42 | 0.00 | 5.82 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 65.75 |
| 06:24:37 PM | 35 | 28.87 | 0.00 | 5.50 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 65.64 |
| 06:24:37 PM | 36 | 29.01 | 0.00 | 6.48 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 64.51 |
| 06:24:37 PM | 37 | 29.15 | 0.00 | 5.76 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 65.08 |
| 06:24:37 PM | 38 | 27.65 | 0.00 | 6.14 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 66.21 |
| 06:24:37 PM | 39 | 27.74 | 0.00 | 6.16 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 66.10 |
| 06:24:37 PM | 40 | 30.17 | 0.00 | 5.42 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 64.41 |
| 06:24:37 PM | 41 | 28.77 | 0.00 | 5.82 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 65.41 |
| 06:24:37 PM | 42 | 28.57 | 0.00 | 6.12 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 65.31 |
| 06:24:37 PM | 43 | 29.59 | 0.00 | 5.78 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 64.63 |
| 06:24:37 PM | 44 | 28.72 | 0.00 | 5.54 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 65.74 |
| 06:24:37 PM | 45 | 28.77 | 0.00 | 5.48 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 65.75 |
| 06:24:37 PM | 46 | 27.68 | 0.00 | 6.57 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 65.74 |
| 06:24:37 PM | 47 | 28.91 | 0.00 | 7.82 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 63.27 |
| 06:24:37 PM | 48 | 75.93 | 0.00 | 11.53 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 12.54 |
| 06:24:37 PM | 49 | 76.53 | 0.00 | 11.56 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 11.90 |
| 06:24:37 PM | 50 | 76.01 | 0.00 | 11.15 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 12.84 |
| 06:24:37 PM | 51 | 76.19 | 0.00 | 11.22 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 12.59 |
| 06:24:37 PM | 52 | 74.74 | 0.00 | 12.63 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 12.63 |
| 06:24:37 PM | 53 | 75.43 | 0.00 | 11.26 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 13.31 |
| 06:24:37 PM | 54 | 73.81 | 0.00 | 12.93 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 13.27 |
| 06:24:37 PM | 55 | 75.09 | 0.00 | 11.95 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 12.97 |
| 06:24:37 PM | 56 | 75.26 | 0.00 | 11.34 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 13.40 |
| 06:24:37 PM | 57 | 74.74 | 0.00 | 11.60 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 13.65 |
| 06:24:37 PM | 58 | 73.97 | 0.00 | 12.67 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 13.36 |
| 06:24:37 PM | 59 | 73.56 | 0.00 | 12.54 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 13.90 |
| 06:24:37 PM | 60 | 75.51 | 0.00 | 11.22 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 13.27 |
| 06:24:37 PM | 61 | 74.32 | 0.00 | 11.99 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 13.70 |
| 06:24:37 PM | 62 | 75.25 | 0.00 | 10.85 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 13.90 |
| 06:24:37 PM | 63 | 72.16 | 0.00 | 10.65 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 17.18 |
| 06:24:37 PM | 64 | 70.21 | 0.00 | 11.64 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 18.15 |
| 06:24:37 PM | 65 | 71.33 | 0.00 | 10.24 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 18.43 |
| 06:24:37 PM | 66 | 71.58 | 0.00 | 10.27 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 18.15 |
| 06:24:37 PM | 67 | 70.21 | 0.00 | 10.62 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 19.18 |
| 06:24:37 PM | 68 | 70.55 | 0.00 | 10.62 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 18.84 |
| 06:24:37 PM | 69 | 70.00 | 0.00 | 11.03 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 18.97 |
| 06:24:37 PM | 70 | 70.21 | 0.00 | 10.96 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 18.84 |
| 06:24:37 PM | 71 | 71.43 | 0.00 | 10.88 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 17.69 |
| 06:24:37 PM | 72 | 31.74 | 0.00 | 8.19 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 60.07 |
| 06:24:37 PM | 73 | 27.74 | 0.00 | 6.51 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 65.75 |
| 06:24:37 PM | 74 | 30.17 | 0.00 | 8.47 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 61.36 |
| 06:24:37 PM | 75 | 29.59 | 0.00 | 6.12 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 64.29 |
| 06:24:37 PM | 76 | 29.59 | 0.00 | 6.12 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 64.29 |
| 06:24:37 PM | 77 | 29.15 | 0.00 | 6.10 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 64.75 |
| 06:24:37 PM | 78 | 27.68 | 0.00 | 6.57 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 65.74 |
| 06:24:37 PM | 79 | 28.33 | 0.00 | 6.48 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 65.19 |
| 06:24:37 PM | 80 | 28.62 | 0.00 | 9.66 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 61.72 |
| 06:24:37 PM | 81 | 29.21 | 0.00 | 5.50 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 65.29 |
| 06:24:37 PM | 82 | 28.03 | 0.00 | 5.54 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 66.44 |
| 06:24:37 PM | 83 | 28.91 | 0.00 | 5.78 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 65.31 |
| 06:24:37 PM | 84 | 28.52 | 0.00 | 5.84 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 65.64 |
| 06:24:37 PM | 85 | 28.91 | 0.00 | 6.12 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 64.97 |
| 06:24:37 PM | 86 | 28.67 | 0.00 | 6.14 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 65.19 |
| 06:24:37 PM | 87 | 28.18 | 0.00 | 5.84 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 65.98 |
| 06:24:37 PM | 88 | 27.78 | 0.00 | 5.90 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 66.32 |
| 06:24:37 PM | 89 | 29.45 | 0.00 | 5.48 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 65.07 |
| 06:24:37 PM | 90 | 28.03 | 0.00 | 6.92 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 65.05 |
| 06:24:37 PM | 91 | 28.18 | 0.00 | 6.53 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 65.29 |
| 06:24:37 PM | 92 | 28.87 | 0.00 | 5.84 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 65.29 |
| 06:24:37 PM | 93 | 28.62 | 0.00 | 5.17 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 66.21 |
| 06:24:37 PM | 94 | 28.91 | 0.00 | 6.12 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 64.97 |
| 06:24:37 PM | 95 | 29.79 | 0.00 | 5.14 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 65.07 |

# Cause

CFS wakeups actively pull wakee tasks
➢ Frequent wakeups from network ISR
➢ Due to the network IRQ binding, waking CPUs are mostly the ones network IRQs are bound to.
➢ Work against periodic and idle load balancing

select_task_rq_fair() has a two-pass process determining whether to wake affine or not
➢ wake_wide() is the first pass, a heuristic that makes sense if waker and wakee are related.
➢ In our cases, waker task is not the one wakes the wakee. It's just happened to be on the CPU when the interrupt comes in.
➢ Wake_wide() returns 0 because waker and wakee have similar wakee_flips numbers.
➢ We notice wake_wide() is the more dominate factor than the second pass, wake_affine()

# Fix?

Questions to be answered:

1. When should we pull for interrupts?
   - Ultimately who has the warmer cache? The scheduler currently doesn't have the necessary information to make a good decision
   - Can we allow the userspace to have a preference?
2. Currently wake_wide() doesn't make sense for wakeups from ISRs. Can we have a better heuristic for interrupts?