1      # Kullback Leibler Divergence in Complete Bacterial and Phage Genomes

2      **Sajia Akhter, Ramy K. Aziz, Mona T. Kashef, Eslam S. Ibrahim, Barbara Bailey, Robert A. Edwards**

3

4      **Supplemental Material**

5

6      **T-test for testing the relationship between GC content and amino acid variations for bacteria and phage samples:**

7

8      There are two samples: bacteria and phages.

9      The sample size of bacteria, $m = 372$

10     The linear model of bacterial sample: $y = 2x^2 - 2x + 0.5$          …     (1)

11     The sample size of phages, $n = 835$

12     The linear model of phage sample: $y = 1.7x^2 - 1.7x + 0.44$       …     (2)

13

14     The general form of these equations is: $y = \beta_0 + \beta_1 x + \beta_2 x^2$

15     We want to test whether the coefficients ($\beta_1$ and $\beta_2$) of these two equations are significantly different:

16

17     $H_0: \beta_1^{\ 1} = \beta_1^{\ 2}$ and $\beta_2^{\ 1} = \beta_2^{\ 2}$

18     $H_A: \beta_1^{\ 1} \neq \beta_1^{\ 2}$ and $\beta_2^{\ 1} \neq \beta_2^{\ 2}$

19     where $\beta_1^{\ 1}$ is the $\beta_1$ coefficient of equation 1 (bacterial sample) and $\beta_1^{\ 2}$ is the $\beta_1$ coefficient of equation 2 (phage sample).

20

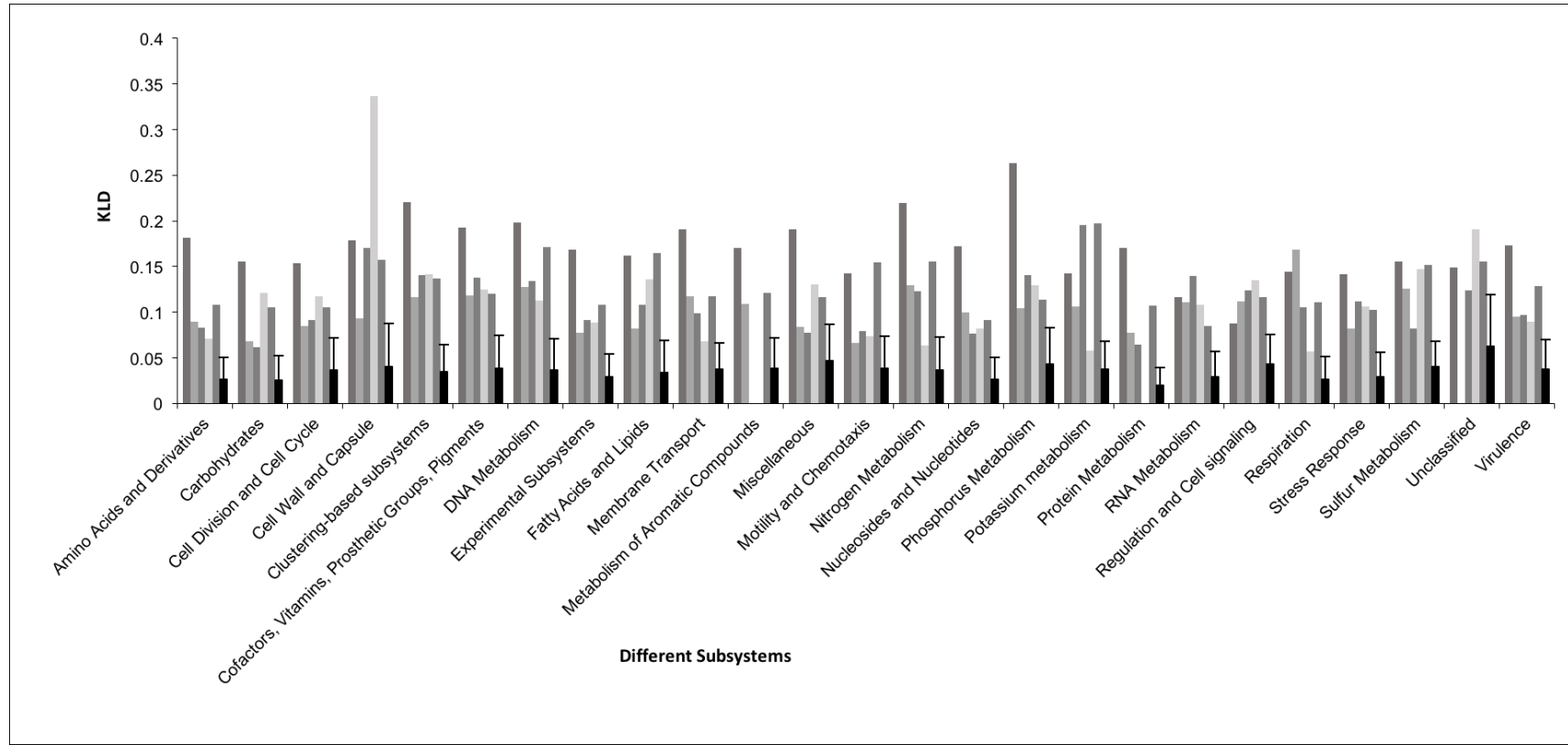21     T test for two independent unequal sample sizes:

22     $$t = \frac{\bar{\beta_p} - \bar{\beta_b}}{\sqrt{SE(\beta_p)^2 + SE(\beta_b)^2}}$$ , where SE is the standard error.

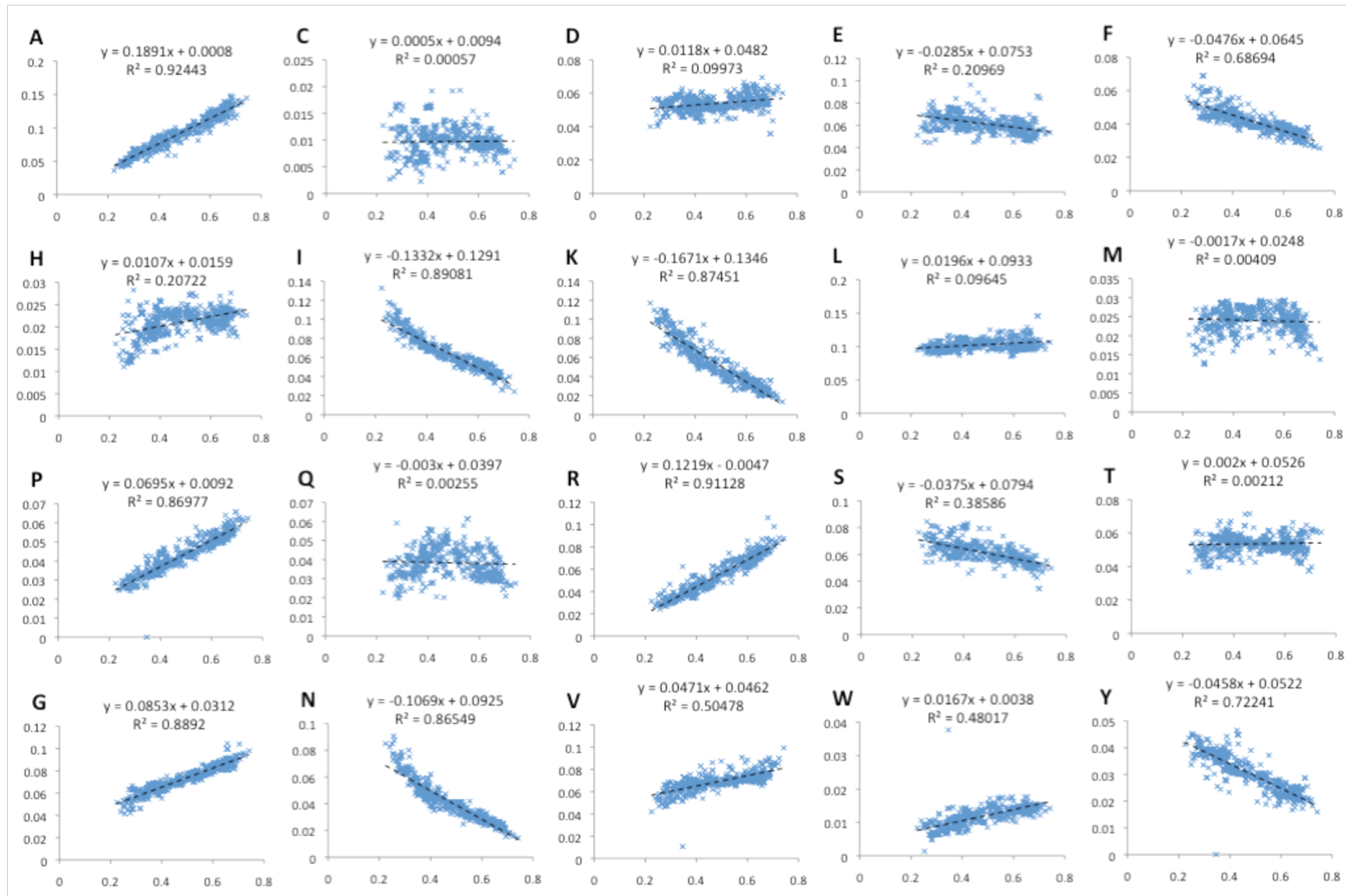23     Here, degree of freedom, $df = m-2 + n-2 = 372-3 + 835-3 = 1201$

24

25     For coefficient $\beta_1$, p value is $1.058481e-07$, and for coefficient $\beta_2$, p value is $1.631291e-06$. So we can reject the null hypothesis.

26

27     This means Equation 1 and Equation 2 are significantly different.

28    **Supplemental Figure 1: Comparison of the divergence of amino acid composition and the phylogenetic group for the most**
29    **divergent bacterial genomes.** The first five bars represent *Wigglesworthia glossinidia*, *Borrelia garinii*, *Mycoplasma mycoides*,
30    *Ureaplasma parvum* serovar and *Buchnera aphidicola*. The sixth bar is for the mean of amino acid utilization for each subsystem.
31
32



33
34
35
36

37    **Supplemental Figure 2: Amino acid frequency of 446 bacterial genomes vs. GC%**
38
39

40    **Supplemental Figure 3: Amino acid frequency of 835 complete phage genomes vs. GC%**

41