

Supplementary material

“BiTNet: Deep Hybrid Model for Ultrasonography Image Analysis of Human Biliary Tract and Its Applications”

Thanapong Intharah, Yupaporn Wanna, Kannika Wiratchawa, Prem Junsawang, Attapol Titapun, Anchalee Techasen, Arunnit Boonrod, Vallop Laopaiboon, Nittaya Chamadol, and Narong Khuntikeo

I. COMPARES THE PERFORMANCE BETWEEN EFFICIENTNET BASE MODEL AND BITNET MODEL MODIFICATION ON 8-FOLD CROSS-VALIDATION SET

A. Compares the median of accuracy between *EfficientNet* model and *BiTNeT* model

1) Null and Alternative Hypotheses

$$H_0 : \theta_1 = \theta_2$$

$$H_1 : \theta_1 \neq \theta_2$$

Where

θ_1 = Median of accuracy of *EfficientNet* model.

θ_2 = Median of accuracy of *BiTNeT* model.

2) The Assumption tests

- There is no relationship of accuracy between *EfficientNet* model and *BiTNeT* model.
- Test of Normality: We use **Shapiro-wilk test** to testing normal distribution of accuracy score each model.

EfficientNet model:

Hypothesis:

H_0 : Accuracy scores of *EfficientNet* model follows normal distribution.

H_1 : Accuracy scores of *EfficientNet* does not follows normal distribution.

Table 1: Result of Test of Normality of accuracy scores of *EfficientNet* model.

	Shapiro-wilk	
	W-test statistic	P-value
EfficientNet model	0.86	0.12
* 99.90% confidence intervals (99.90% CI) and p-values from testing ($p \leq 0.001$ was considered statistically significant).		

The test is non-significant, $W = 0.860$, $p = 0.120$, which indicates that the accuracy scores of *EfficientNet* model are normally distributed.

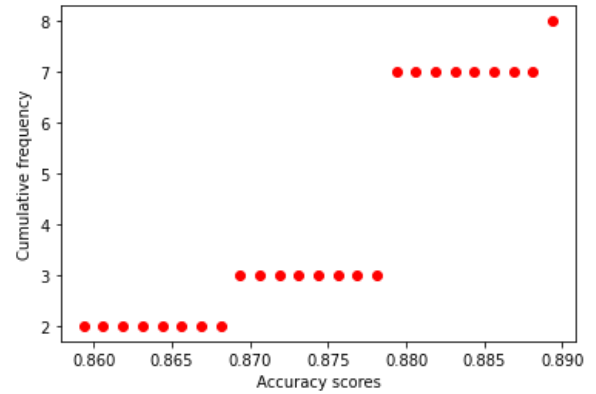


Fig 1: Probability Plots (PP Plot) of accuracy scores of *EfficientNet* model.

BiTNeT model:

Hypothesis:

H_0 : Accuracy scores of *BiTNeT* model follows normal distribution.

H_1 : Accuracy scores of *BiTNeT* model follows normal distribution does not follows normal distribution.

Table 2: Result of Test of Normality of accuracy scores of *EfficientNet* model.

	Shapiro-wilk	
	W-test statistic	P-value
BiTNeT model	0.66	0.00
* 99.90% confidence intervals (99.90% CI) and p-values from testing ($p \leq 0.001$ was considered statistically significant).		

The test had a significant, $W = 0.665$, $p = 0.000$, which indicates that the accuracy scores of *BiTNeT* model does not follows normal distribution.

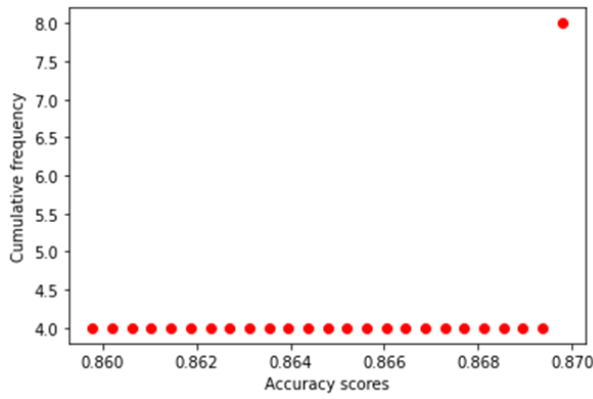


Fig 2: Probability Plots (PP Plot) of accuracy scores of BiTNet model.

3) Test Statistics

The test statistic for this the **Mann Whitney U-Test** is denoted U, for compare group rank differences.

Table 3: Result of Mann Whitney U-Test between EfficientNet model and BiTNet model: accuracy scores.

Mann-Whitney Test	
	EfficientNet model × BiTNet model
Mann Whitney U	50.00
P-value	0.05
*With 99.90% confidence intervals (99.90% CI) and p-values from testing (a two - tailed $p \leq 0.001$ was considered statistically significant).	

B. Compares the mean of precision between EfficientNet model and BiTNet model

1) Null and Alternative Hypotheses

$$H_0 : \mu_1 = \mu_2$$

$$H_1 : \mu_1 \neq \mu_2$$

Where

μ_1 = Mean of precision of EfficientNet model.

μ_2 = Mean of precision of BiTNet model.

2) The Assumption tests

- There is no relationship of precision between EfficientNet model and BiTNet model.
- Test of Normality: We use **Shapiro-wilk test** to testing normal distribution of precision scores each model.

EfficientNet model:

Hypothesis:

H_0 : Precision scores of EfficientNet model follows normal distribution.

H_1 : Precision scores of EfficientNet does not follows normal distribution.

Table 4: Result of Test of Normality of precision scores of EfficientNet model.

	Shapiro-wilk	
	W-test statistic	P-value
EfficientNet model	0.89	0.23
* 99.90% confidence intervals (99.90% CI) and p-values from testing ($p \leq 0.001$ was considered statistically significant).		

The test is non-significant, $W = 0.89$, $p = 0.23$, which indicates that the precision scores of EfficientNet model follow normally distributed.

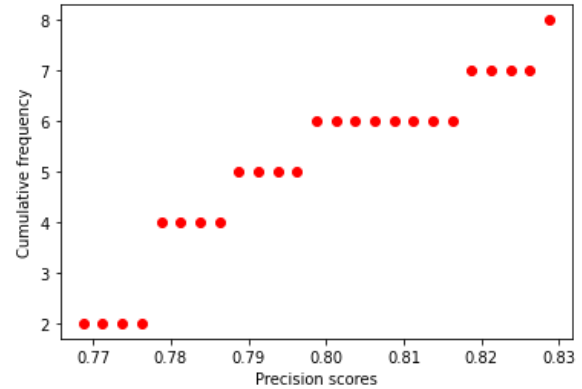


Fig 3: Probability Plots (PP Plot) of precision scores of EfficientNet model.

BiTNet model:

Hypothesis:

H_0 : Precision scores of BiTNet model follows normal distribution.

H_1 : Precision scores of BiTNet does not follows normal distribution.

Table 4: Result of Test of Normality of precision scores of BiTNet model.

	Shapiro-wilk	
	W-test statistic	P-value
BiTNet model	0.88	0.21
* 99.90% confidence intervals (99.90% CI) and p-values from testing ($p \leq 0.001$ was considered statistically significant).		

The test is non-significant, $W = 0.88$, $p = 0.21$, which indicates that the precision scores of BiTNet model follow normally distributed.

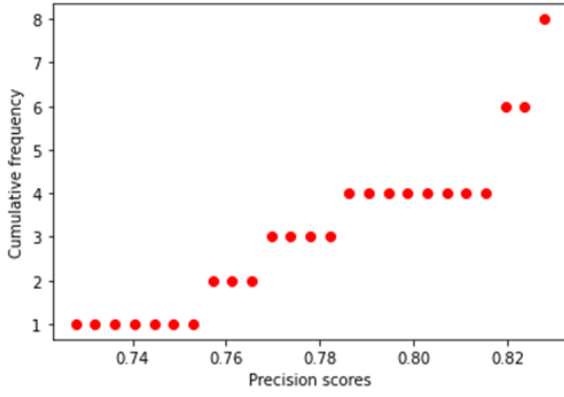


Fig 3: Probability Plots (PP Plot) of precision scores of BiTNet model.

- Test of Homogeneity of variances: We use **Levene's Test** to test for the homogeneity of variance of the precision between EfficientNet model and BiTNet model.

Hypothesis

$$H_0 : \sigma_1^2 - \sigma_2^2 = 0$$

$$H_1 : \sigma_1^2 - \sigma_2^2 \neq 0$$

Where

σ_1^2 = Variances of the precision of EfficientNet model.

σ_2^2 = Variances of the precision of BiTNet model.

Table 5: Result of Test for Equality of Variances of precision between EfficientNet model and BiTNet model.

	Levene's Test for Equality of Variances	
	F	P-value
Equal variance assumed	3.33	0.08
* 99.90% confidence intervals (99.90% CI) and p-values from testing ($p \leq 0.001$ was considered statistically significant).		

The test is non-significant, $F = 3.33$, $p = 0.08$, which indicates that the population variances of precision between EfficientNet model and BiTNet model are equal. When equal variances are assumed, the calculation uses pooled variances to use Independent Samples T-Test

3) Test Statistics

The test statistic for this **Independent Samples T-Test** is denoted t , for equal variances are assumed.

Table 6: Result of Independent Samples T-Test between EfficientNet model and BiTNet model: precision scores.

Two sample t-test with equal variance				
P - value	t	Mean difference	99.90% Confident Interval of the difference	
			Lower	Upper
0.94	-0.08	-1.25×10^{-3}	-0.06	0.06
*With 99.90% confidence intervals (99.90% CI) and p-values from testing (a two-tailed $p \leq 0.001$ was considered statistically significant).				

4) Interval estimates Using T-score with 99.90% CI

Table 7: Result of Interval estimates of precision scores using T-score.

Interval estimates using T-score			
Model	Mean of precision scores	99.90% Confident Interval	
		Lower	Upper
EfficientNet	79.25	74.41	84.09
BiTNet	79.37	71.33	87.42

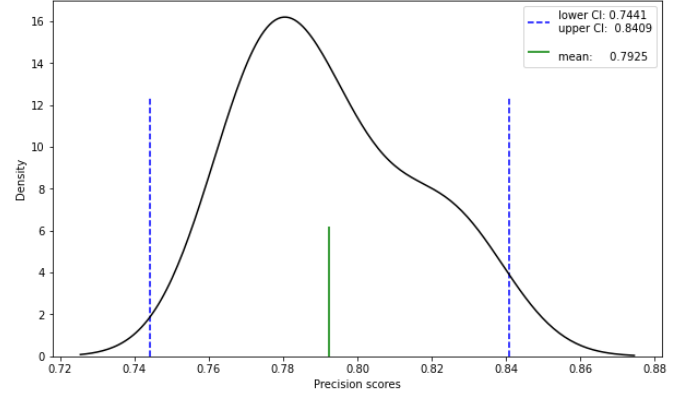


Fig 4: Plot of precision scores of EfficientNet model, t-statistics - Confidence Level = 99.90%.

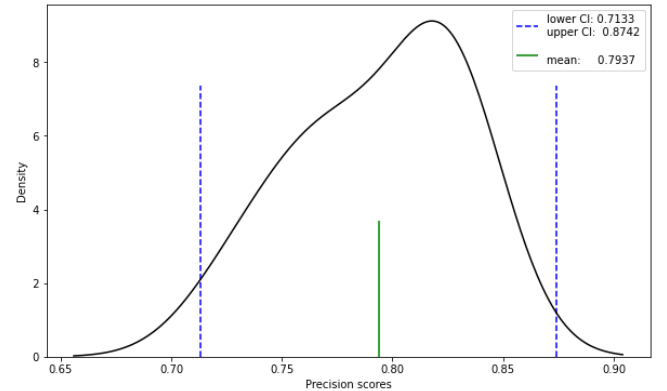


Fig 5: Plot of precision scores of BiTNet model, t-statistics - Confidence Level = 99.90%.

C. Compares the mean of recall between EfficientNet model and BiTNet model

1) Null and Alternative Hypotheses

$$H_0 : \mu_1 = \mu_2$$

$$H_1 : \mu_1 \neq \mu_2$$

Where

μ_1 = Mean of recall of EfficientNet model.

μ_2 = Mean of recall of BiTNet model.

2) The Assumption tests

- There is no relationship of recall between EfficientNet model and BiTNet model.

- Test of Normality: We use **Shapiro-wilk test** to testing normal distribution of recall scores each model.

EfficientNet model:

Hypothesis:

H_0 : Recall scores of EfficientNet model follows normal distribution.

H_1 : Recall scores of EfficientNet does not follows normal distribution.

Table 8: Result of Test of Normality of recall scores of EfficientNet model.

	Shapiro-wilk	
	W-test statistic	P-value
EfficientNet model	0.96	0.85
* 99.90% confidence intervals (99.90% CI) and p-values from testing ($p \leq 0.001$ was considered statistically significant).		

The test is non-significant, $W = 0.96$, $p = 0.85$, which indicates that the recall scores of EfficientNet model follow normally distributed.

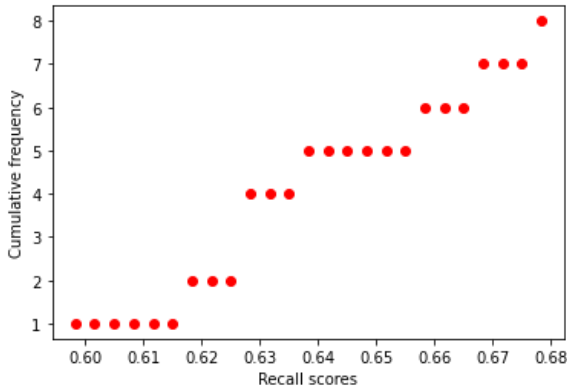


Fig 6: Probability Plots (PP Plot) of recall scores of EfficientNet model.

BiTNet model:

Hypothesis:

H_0 : Recall scores of BiTNet model follows normal distribution.

H_1 : Recall scores of BiTNet does not follows normal distribution.

Table 9: Result of Test of Normality of recall scores of BiTNet model.

	Shapiro-wilk	
	W-test statistic	P-value
BiTNet model	0.97	0.93
* 99.90% confidence intervals (99.90% CI) and p-values from testing ($p \leq 0.001$ was considered statistically significant).		

The test is non-significant, $W = 0.97$, $p = 0.93$, which indicates that the recall scores of BiTNet model follow normally distributed.

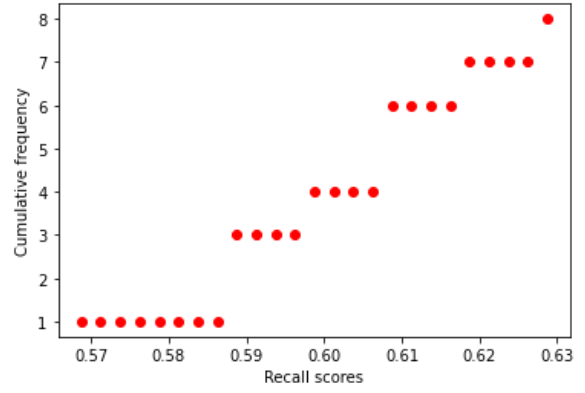


Fig 7: Probability Plots (PP Plot) of recall scores of BiTNet model.

- Test of Homogeneity of variances: We use **Levene's Test** to test for the homogeneity of variance of the recall between EfficientNet model and BiTNet model.

Hypothesis

$$H_0 : \sigma_1^2 - \sigma_2^2 = 0$$

$$H_1 : \sigma_1^2 - \sigma_2^2 \neq 0$$

Where

σ_1^2 = Variances of the recall of EfficientNet model.

σ_2^2 = Variances of the recall of BiTNet model.

Table 10: Result of Test for Equality of Variances of recall between EfficientNet model and BiTNet model.

	Levene's Test for Equality of Variances	
	F	P-value
Equal variance assumed	1.14	0.30
* 99.90% confidence intervals (99.90% CI) and p-values from testing ($p \leq 0.001$ was considered statistically significant).		

The test is non-significant, $F = 1.14$, $p = 0.30$, which indicates that the population variances of recall between EfficientNet model and BiTNet model are equal. When equal variances are assumed, the calculation uses pooled variances to use Independent Samples T-Test

3) Test Statistics

The test statistic for this **Independent Samples T-Test** is denoted t , for equal variances are assumed.

Table 11: Result of Independent Samples T-Test between EfficientNet model and BiTNet model: recall scores

Two sample t-test with equal variance				
P - value	t	Mean difference	99.90% Confident Interval of the difference	
			Lower	Upper
5.07×10^{-3}	3.32	0.04	-9.60×10^{-3}	0.09
*With 99.90% confidence intervals (99.90% CI) and p-values from testing (a two-tailed $p \leq 0.001$ was considered statistically significant).				

4) Interval estimates Using T-score with 99.90% CI

Table 12: Result of Interval estimates of recall scores using T-score.

Interval estimates using T-score			
Model	Mean of recall scores	99.90% Confident Interval	
		Lower	Upper
EfficientNet	0.64	0.60	0.70
BiTNet	0.60	0.56	0.64

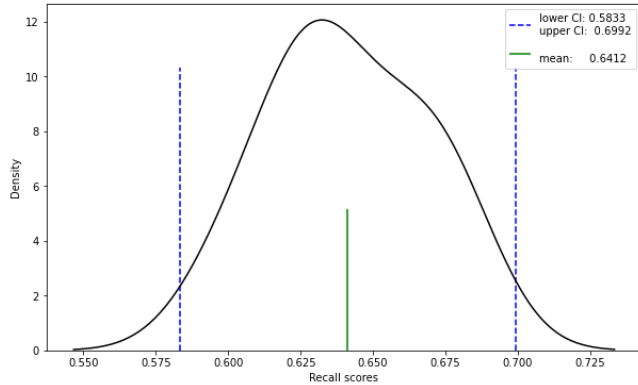


Fig 8: Plot of recall scores of EfficientNet model, t-statistics - Confidence Level = 99.90%.

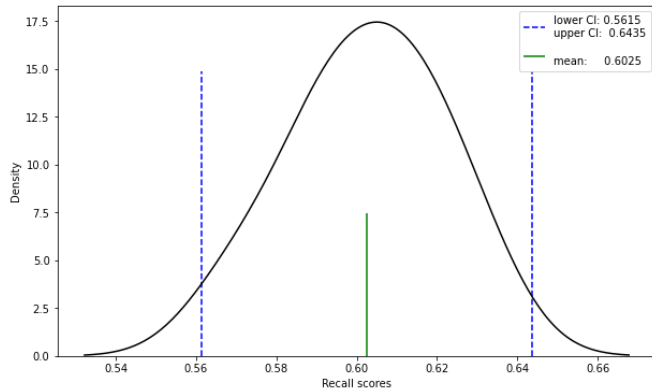


Fig 9: Plot of precision scores of BiTNet model, t-statistics - Confidence Level = 99.90%.

II. COMPARES THE PERFORMANCE BETWEEN EFFICIENTNET BASE MODEL AND BITNET MODEL MODIFICATION ON 8-FOLD TEST SET

A. Compares the mean of accuracy between EfficientNet model and BiTNet model

1) Null and Alternative Hypotheses

$$H_0 : \mu_1 = \mu_2$$

$$H_1 : \mu_1 \neq \mu_2$$

Where

μ_1 = Mean of accuracy of EfficientNet model.

μ_2 = Mean of accuracy of BiTNet model.

2) The Assumption tests

- There is no relationship of accuracy between EfficientNet model and BiTNet model.
- Test of Normality: We use **Shapiro-wilk test** to testing normal distribution of accuracy scores each model.

EfficientNet model:

Hypothesis:

H_0 : Accuracy scores of EfficientNet model follows normal distribution.

H_1 : Accuracy scores of EfficientNet does not follows normal distribution.

Table 13: Result of Test of Normality of accuracy scores of EfficientNet model.

	Shapiro-wilk	
	W-test statistic	P-value
EfficientNet model	0.83	0.05
* 99.90% confidence intervals (99.90% CI) and p-values from testing ($p \leq 0.001$ was considered statistically significant).		

The test is non-significant, $W = 0.83$, $p = 0.05$, which indicates that the accuracy scores of EfficientNet model follow normally distributed.

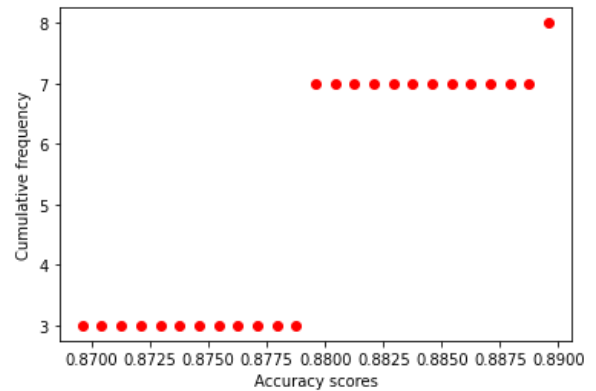


Fig 10: Probability Plots (PP Plot) of accuracy scores of EfficientNet model.

BiTNet model:

Hypothesis:

H_0 : Accuracy scores of BiTNet model follows normal distribution.

H_1 : Accuracy scores of BiTNet does not follows normal distribution.

Table 14: Result of Test of Normality of precision scores of BiTNet model.

	Shapiro-wilk	
	W-test statistic	P-value
BiTNet model	0.80	0.03
* 99.90% confidence intervals (99.90% CI) and p-values from testing ($p \leq 0.001$ was considered statistically significant).		

The test is non-significant, $W = 0.80$, $p = 0.03$, which

indicates that the accuracy scores of BiTNet model follow normally distributed.

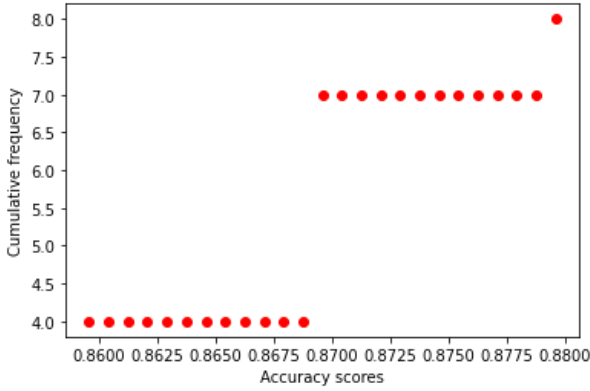


Fig 11: Probability Plots (PP Plot) of accuracy scores of BiTNet model.

- Test of Homogeneity of variances: We use **Levene's Test** to test for the homogeneity of variance of the accuracy between EfficientNet model and BiTNet model.

Hypothesis

$$H_0 : \sigma_1^2 - \sigma_2^2 = 0$$

$$H_1 : \sigma_1^2 - \sigma_2^2 \neq 0$$

Where

σ_1^2 = Variances of the accuracy of EfficientNet model.

σ_2^2 = Variances of the accuracy of BiTNet model.

Table 15: Result of Test for Equality of Variances of accuracy between EfficientNet model and BiTNet model.

	Levene's Test for Equality of Variances	
	F	P-value
Equal variance assumed	0.13	0.73
* 99.90% confidence intervals (99.90% CI) and p-values from testing ($p \leq 0.001$ was considered statistically significant).		

The test is non-significant, $F = 0.13$, $p = 0.73$, which indicates that the population variances of accuracy between EfficientNet and BiTNet model are equal. When equal variances are assumed, the calculation uses pooled variances to use Independent Samples T-Test

3) Test Statistics

The test statistic for this **Independent Samples T-Test** is denoted t , for equal variances are assumed.

Table 16: Result of Independent Samples T-Test between EfficientNet model and BiTNet model: accuracy scores.

Two sample t-test with equal variance				
P - value	t	Mean difference	99.90% Confident Interval of the difference	
			Lower	Upper
0.01	3.10	0.01	-3.77×10^{-3}	0.03
*With 99.90% confidence intervals (99.90% CI) and p-values from testing (a two-tailed $p \leq 0.001$ was considered statistically significant).				

4) Interval estimates Using T-score with 99.90% CI

Table 17: Result of Interval estimates of accuracy scores using T-score.

Interval estimates using T-score			
Model	Mean of accuracy scores	99.90% Confident Interval	
		Lower	Upper
EfficientNet	87.75	86.23	89.27
BiTNet	86.63	85.03	88.22

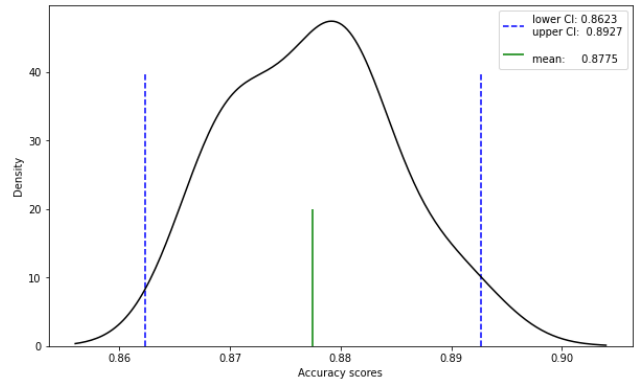


Fig 12: Plot of accuracy scores of EfficientNet model, t-statistics - Confidence Level = 99.90%.

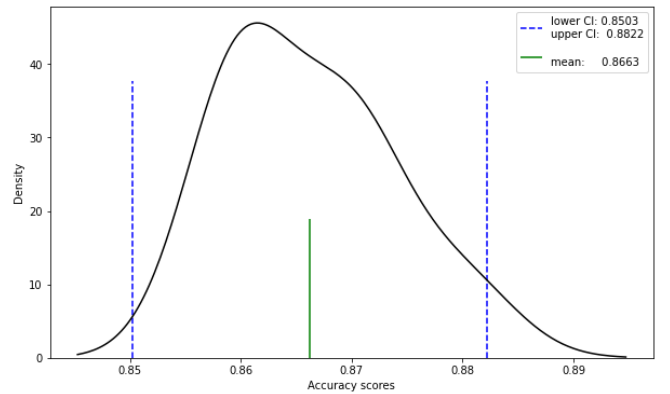


Fig 13: Plot of accuracy scores of BiTNet model, t-statistics - Confidence Level = 99.90%.

B. Compares the mean of precision between EfficientNet model and BiTNet model

1) Null and Alternative Hypotheses

$$H_0 : \mu_1 = \mu_2$$

$$H_1 : \mu_1 \neq \mu_2$$

Where

μ_1 = Mean of precision of EfficientNet model.

μ_2 = Mean of precision of BiTNet model.

2) The Assumption tests

- There is no relationship of precision between EfficientNet model and BiTNet model.
- Test of Normality: We use **Shapiro-wilk test** to testing normal distribution of precision scores each model.

EfficientNet model:

Hypothesis:

H_0 : Precision scores of EfficientNet model follows normal distribution.

H_1 : Precision scores of EfficientNet does not follows normal distribution.

Table 18: Result of Test of Normality of precision scores of EfficientNet model.

	Shapiro-wilk	
	W-test statistic	P-value
EfficientNet model	0.87	0.15
* 99.90% confidence intervals (99.90% CI) and p-values from testing ($p \leq 0.001$ was considered statistically significant).		

The test is non-significant, $W = 0.87$, $p = 0.15$, which indicates that the precision scores of EfficientNet model follow normally distributed.

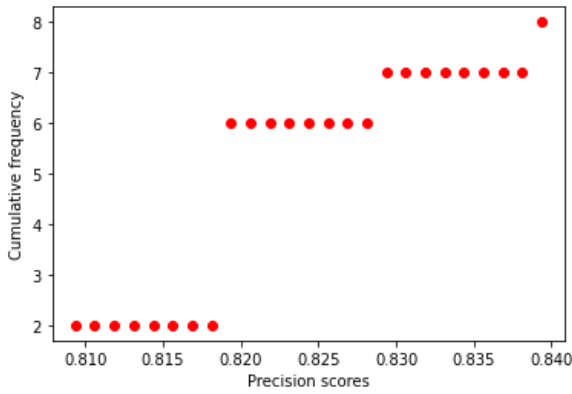


Fig 14: Probability Plots (PP Plot) of precision scores of EfficientNet model.

BiTNet model:

Hypothesis:

H_0 : Precision scores of BiTNet model follows normal distribution.

H_1 : Precision scores of BiTNet does not follows normal distribution.

Table 19: Result of Test of Normality of precision scores of BiTNet model.

	Shapiro-wilk	
	W-test statistic	P-value
BiTNet model	0.87	0.15
* 99.90% confidence intervals (99.90% CI) and p-values from testing ($p \leq 0.001$ was considered statistically significant).		

The test is non-significant, $W = 0.87$, $p = 0.15$, which indicates that the precision scores of BiTNet model follow normally distributed.

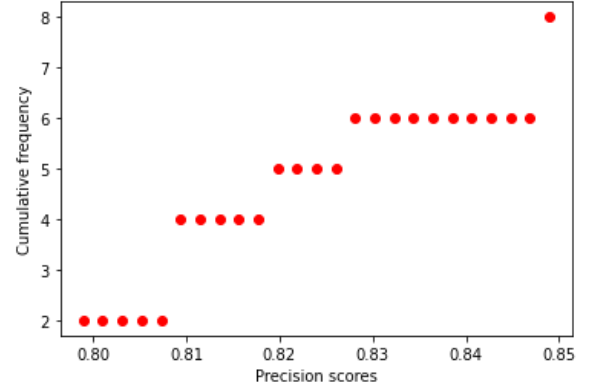


Fig 15: Probability Plots (PP Plot) of precision scores of BiTNet model.

- Test of Homogeneity of variances: We use **Levene's Test** to test for the homogeneity of variance of the precision between EfficientNet model and BiTNet model.

Hypothesis

$$H_0 : \sigma_1^2 - \sigma_2^2 = 0$$

$$H_1 : \sigma_1^2 - \sigma_2^2 \neq 0$$

Where

σ_1^2 = Variances of the precision of EfficientNet model.

σ_2^2 = Variances of the precision of BiTNet model.

Table 20: Result of Test for Equality of Variances of precision between EfficientNet model and BiTNet model.

	Levene's Test for Equality of Variances	
	F	P-value
Equal variance assumed	5.24	0.04
* 99.90% confidence intervals (99.90% CI) and p-values from testing ($p \leq 0.001$ was considered statistically significant).		

The test is non-significant, $F = 5.24$, $p = 0.04$, which indicates that the population variances of precision between EfficientNet model and BiTNet model are equal. When equal variances are assumed, the calculation uses pooled variances to use Independent Samples T-Test

3) Test Statistics

The test statistic for this **Independent Samples T-Test** is denoted t , for equal variances are assumed.

Table 21: Result of Independent Samples T-Test between EfficientNet model and BiTNet model: precision scores.

Two sample t-test with equal variance				
P - value	t	Mean difference	99.90% Confident Interval of the difference	
			Lower	Upper
1.00	0.00	0.00	-0.03	0.03
*With 99.90% confidence intervals (99.90% CI) and p-values from testing (a two-tailed $p \leq 0.001$ was considered statistically significant).				

4) Interval estimates Using T-score with 99.90% CI

Table 22: Result of Interval estimates of precision scores using T-score.

Interval estimates using T-score			
Model	Mean of precision scores	99.90% Confident Interval	
		Lower	Upper
EfficientNet	82.13	79.99	84.26
BiTNet	82.13	77.76	86.49

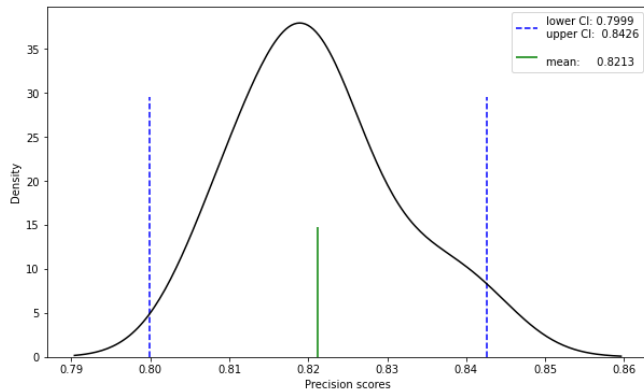


Fig 16: Plot of precision scores of EfficientNet model, t-statistics - Confidence Level = 99.90%.

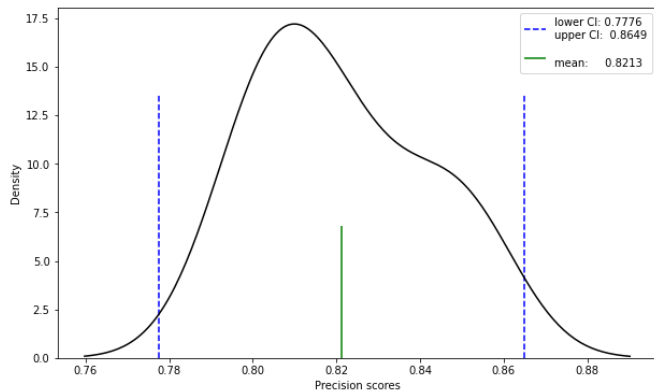


Fig 17: Plot of precision scores of BiTNet model, t-statistics - Confidence Level = 99.90%.

C. Compares the mean of recall between EfficientNet model and BiTNet model

1) Null and Alternative Hypotheses

$$H_0 : \mu_1 = \mu_2$$

$$H_1 : \mu_1 \neq \mu_2$$

Where

μ_1 = Mean of recall of EfficientNet model.

μ_2 = Mean of recall of BiTNet model.

2) The Assumption tests

- There is no relationship of recall between EfficientNet model and BiTNet model.
- Test of Normality: We use **Shapiro-wilk test** to testing normal distribution of recall scores each model.

EfficientNet model:

Hypothesis:

H_0 : Recall scores of EfficientNet model follows normal distribution.

H_1 : Recall scores of EfficientNet does not follows normal distribution.

Table 23: Result of Test of Normality of recall scores of EfficientNet model.

	Shapiro-wilk	
	W-test statistic	P-value
EfficientNet model	0.98	0.96
* 99.90% confidence intervals (99.90% CI) and p-values from testing ($p \leq 0.001$ was considered statistically significant).		

The test is non-significant, $W = 0.98$, $p = 0.96$, which indicates that the recall scores of EfficientNet model follow normally distributed.

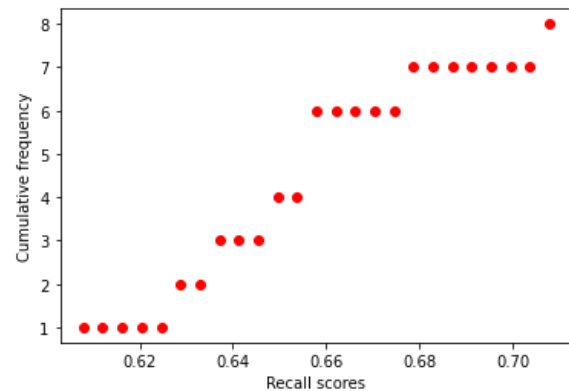


Fig 18: Probability Plots (PP Plot) of recall scores of EfficientNet model.

BiTNet model:

Hypothesis:

H_0 : Recall scores of BiTNet model follows normal distribution.

H_1 : Recall scores of BiTNet does not follows normal distribution.

Table 24: Result of Test of Normality of recall scores of BiTNet model.

	Shapiro-wilk	
	W-test statistic	P-value
BiTNet model	0.95	0.75
* 99.90% confidence intervals (99.90% CI) and p-values from testing ($p \leq 0.001$ was considered statistically significant).		

The test is non-significant, $W = 0.95$, $p = 0.75$, which indicates that the recall scores of BiTNet model follow normally distributed.

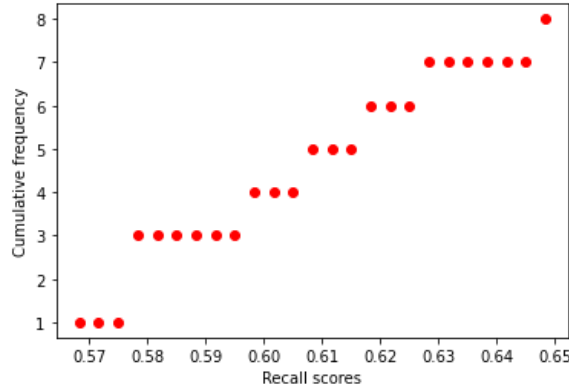


Fig 19: Probability Plots (PP Plot) of recall scores of BiTNet model.

- Test of Homogeneity of variances: We use **Levene's Test** to test for the homogeneity of variance of the recall between EfficientNet model and BiTNet model.

Hypothesis

$$H_0: \sigma_1^2 - \sigma_2^2 = 0$$

$$H_1: \sigma_1^2 - \sigma_2^2 \neq 0$$

Where

σ_1^2 = Variances of the recall of EfficientNet model.

σ_2^2 = Variances of the recall of BiTNet model.

Table 25: Result of Test for Equality of Variances of recall between EfficientNet model and BiTNet model.

	Levene's Test for Equality of Variances	
	F	P-value
Equal variance assumed	0.76×10^{-30}	1.0
* 99.90% confidence intervals (99.90% CI) and p-values from testing ($p \leq 0.001$ was considered statistically significant).		

The test is non-significant, $F = 0.76 \times 10^{-30}$, $p = 1.0$, which indicates that the population variances of recall between EfficientNet model and BiTNet model are equal. When equal variances are assumed, the calculation uses pooled variances to use Independent Samples T-Test

3) Test Statistics

The test statistic for this **Independent Samples T-Test** is denoted t , for equal variances are assumed.

Table 26: Result of Independent Samples T-Test between EfficientNet model and BiTNet model: recall scores

Two sample t-test with equal variance				
P - value	t	Mean difference	99.90% Confident Interval of the difference	
			Lower	Upper
4.20×10^{-3}	3.42	0.05	-0.01	0.11
*With 99.90% confidence intervals (99.90% CI) and p-values from testing (a two-tailed $p \leq 0.001$ was considered statistically significant).				

4) Interval estimates Using T-score with 99.90% CI

Table 27: Result of Interval estimates of recall scores using T-score.

Interval estimates using T-score			
Model	Mean of recall scores	99.90% Confident Interval	
		Lower	Upper
EfficientNet	65.50	58.90	72.10
BiTNet	60.50	54.53	66.47

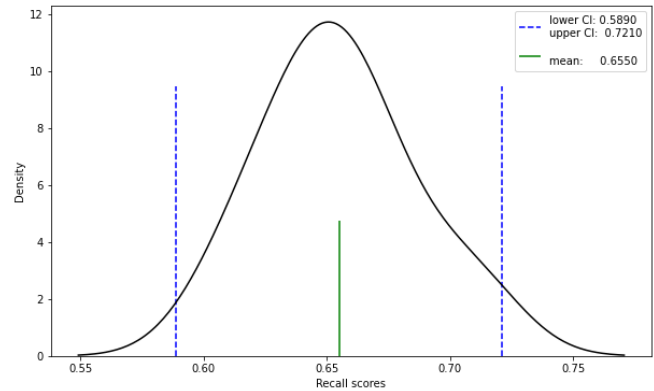


Fig 20: Plot of recall scores of EfficientNet model, t-statistics - Confidence Level = 99.90%.

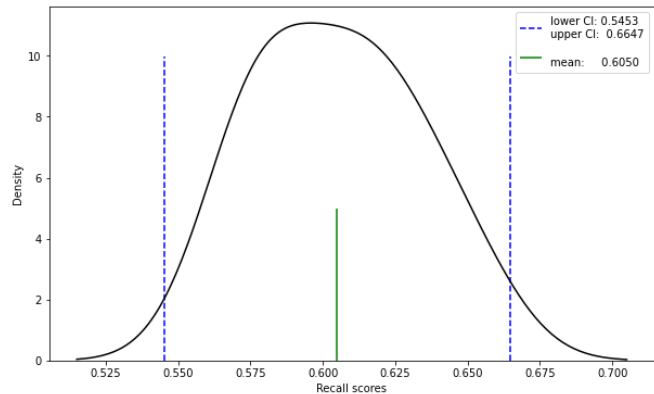


Fig 21: Plot of precision scores of BiTNet model, t-statistics - Confidence Level = 99.90%.

III. COMPARES THE MEAN DIFFERENCES BETWEEN PREDICTION CONFIDENCE OF THE CORRECT AND INCORRECT GROUPS

We use **Independent Samples T-Test** to compare the means of mean difference of prediction confidence of the correct and incorrect groups between BiTNet model and EfficientNet model.

3.1 Null and Alternative Hypotheses

$$H_0 : \mu_1 = \mu_2$$

$$H_1 : \mu_1 > \mu_2$$

Where

μ_1 = Mean of mean difference of prediction confidence of BiTNet model.

μ_2 = Mean of mean difference of prediction confidence of EfficientNet model.

3.2 The Assumption tests

1) There is no relationship between mean differences of BiTNet model and mean differences of EfficientNet model.

2) Test of Normality: We use **Shapiro-wilk test** to testing normal distribution of mean difference of prediction confidence each model.

BiTNet model:

Hypothesis:

H_0 : Mean difference of prediction confidence of BiTNet model follows normal distribution.

H_1 : Mean difference of prediction confidence of BiTNet model does not follow normal distribution.

Table 28: Result of Test of Normality of the mean difference of prediction confidence of BiTNet model.

	Shapiro-wilk	
	W-test statistic	P-value
Mean difference	0.92	2.72×10^{-2}
* 99.90% confidence intervals (99.90% CI) and p-values from testing ($p \leq 0.001$ was considered statistically significant).		

The test is non-significant, $W = 0.92$, $p = 2.72 \times 10^{-2}$, which indicates that the mean difference of prediction confidence of BiTNet model are normally distributed.

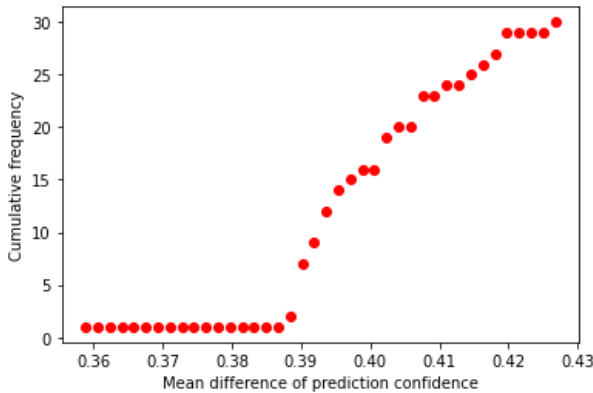


Fig 22: Probability Plots (PP Plot) of the mean difference prediction confidence of BiTNet model.

EfficientNet model:

Hypothesis:

H_0 : Mean difference of prediction confidence of EfficientNet model follows normal distribution.

H_1 : Mean difference of prediction confidence of EfficientNet model does not follow normal distribution.

Table 29: Result of Test of Normality of the mean difference prediction confidence of BiTNet model.

	Shapiro-wilk	
	W-test statistic	P-value
Mean difference	0.93	6.27×10^{-2}
* 99.90% confidence intervals (99.90% CI) and p-values from testing ($p \leq 0.001$ was considered statistically significant).		

The test is non-significant, $W = 0.93$, $p = 6.27 \times 10^{-2}$, which indicates that the mean difference of prediction confidence of EfficientNet model are normally distributed.

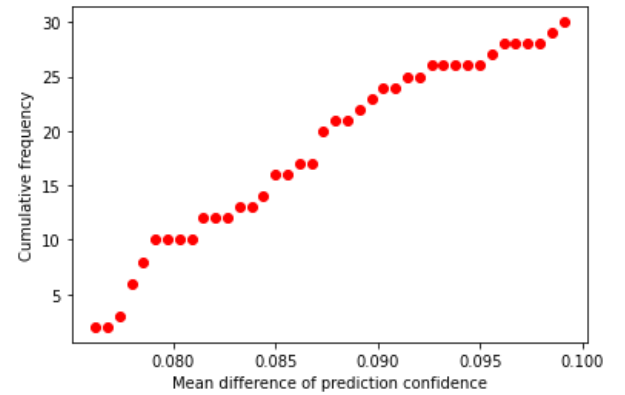


Fig 23: Probability Plots (PP Plot) of the mean difference of prediction confidence of EfficientNet model.

3) Test of Homogeneity of variances

We use **Levene's Test** to test for the homogeneity of variance of the Mean difference of prediction confidence both models

Hypothesis

$$H_0 : \sigma_1^2 - \sigma_2^2 = 0$$

$$H_1 : \sigma_1^2 - \sigma_2^2 \neq 0$$

Where

σ_1^2 = Variances of the Mean difference of prediction confidence of BiTNet model.

σ_2^2 = Variances of the Mean difference of prediction confidence of EfficientNet model.

Table 30: Result of Test for Equality of Variances of the mean difference of prediction confidence between BiTNet model and EfficientNet model.

	Levene's Test for Equality of Variances	
	F	P-value
Equal variance assumed	8.17	5.89×10^{-3}
* 99.90% confidence intervals (99.90% CI) and p-values from testing ($p \leq 0.001$ was considered statistically significant).		

The test is non-significant, $F=8.17$, $p = 5.89 \times 10^{-3}$, which indicates that the population variances of BiTNet model and EfficientNet model are equal. When equal variances are assumed, the calculation uses pooled variances to use Independent Samples T-Test

3.3 Test Statistics

The test statistic for this **Independent Samples T-Test** is denoted t , for equal variances are assumed.

Table 31: Result of Independent Samples T-Test for compare the means of mean difference between BiTNet model and EfficientNet model.

Two sample t-test with equal variance				
P - value	t	Mean difference	99.90% Confident Interval of the difference	
			Lower	Upper
2.34×10^{-70}	114.60	31.58	30.62	32.53
*With 99.90% confidence intervals (99.90% CI) and p-values from testing (a one - tailed $p \leq 0.001$ was considered statistically significant).				

3.4 Interval estimates Using T-score with 99.90% CI

Table 32: Result of Interval estimates of the mean differences using T-score.

Interval estimates using T-score			
Model	Mean of mean difference	99.90% Confident Interval	
		Lower	Upper
BiTNet	40.13	39.29	40.97
EfficientNet	8.55	8.13	8.98

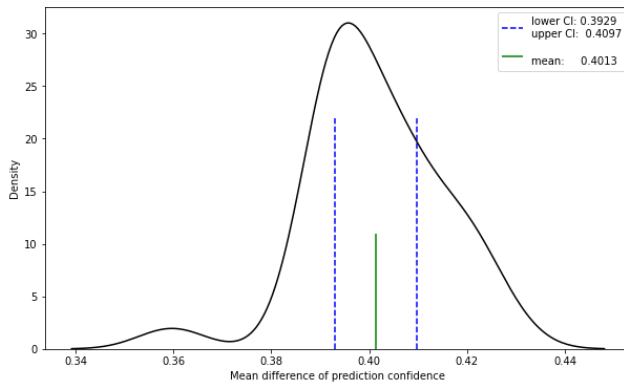


Fig 24: Plot of the mean difference of prediction confidence of the correct and incorrect of BiTNet model, t-statistics - Confidence Level = 99.90%.

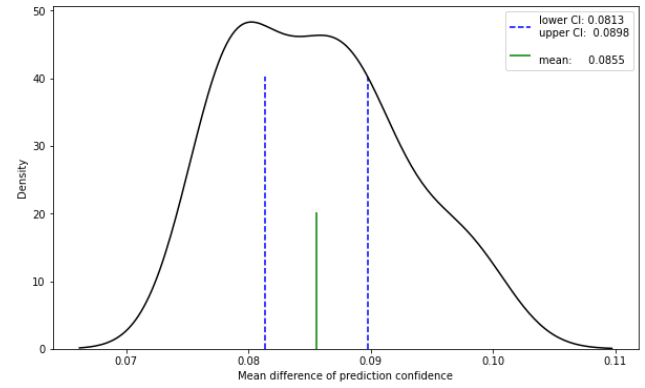


Fig 25: Plot of the mean difference of prediction confidence of the correct and incorrect of EfficientNet model, t-statistics - Confidence Level = 99.90%.

A. Compares the means of prediction confidence between correct and incorrect of BiTNet model

1) Null and Alternative Hypotheses

$$H_0 : \mu_1 = \mu_2$$

$$H_1 : \mu_1 > \mu_2$$

Where

μ_1 = Mean of prediction confidence correct.

μ_2 = Mean of prediction confidence incorrect.

2) The Assumption tests

- There is no relationship of prediction confidence between correct and incorrect.
- Test of Normality: We use **Shapiro-wilk test** to testing normal distribution of mean each prediction confidence.

Prediction confidences correct:

Hypothesis:

H_0 : Mean of prediction confidence correct follows normal distribution.

H_1 : Mean of prediction confidence correct does not follows normal distribution.

Table 33: Result of Test of Normality of prediction confidence correct.

	Shapiro-wilk	
	W-test statistic	P-value
Correct	0.96	0.40
* 99.90% confidence intervals (99.90% CI) and p-values from testing ($p \leq 0.001$ was considered statistically significant).		

The test is non-significant, $W = 0.96$, $p = 0.40$, which indicates that the Mean of confidence correct are normally distributed.

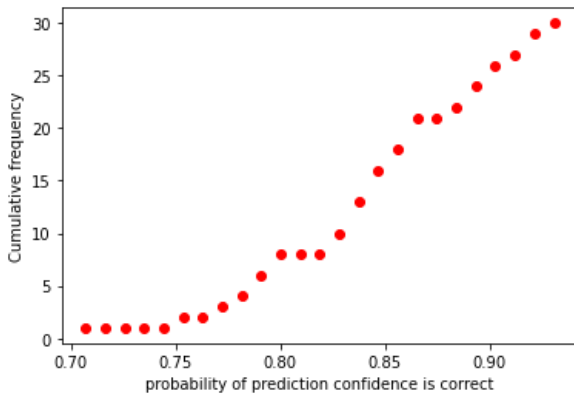


Fig 26: Probability Plots (PP Plot) of prediction confidence is correct.

Prediction confidences incorrect:

Hypothesis:

H_0 : Mean of prediction confidence incorrect follows normal distribution.

H_1 : Mean of prediction confidence incorrect does not follows normal distribution.

Table 34: Result of Test of Normality of prediction confidence incorrect.

	Shapiro-wilk	
	W-test statistic	P-value
Incorrect	0.98	0.72
* 99.90% confidence intervals (99.90% CI) and p-values from testing ($p \leq 0.001$ was considered statistically significant).		

The test is non-significant, $W = 0.98$, $p = 0.72$, which indicates that the mean of confidence incorrect are normally distributed.

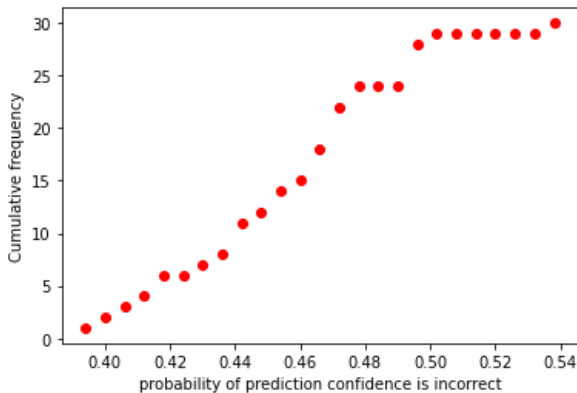


Fig 27: Probability Plots (PP Plot) of prediction confidence incorrect.

- Test of Homogeneity of variances: We use **Levene's Test** to test for the homogeneity of variance of the mean of prediction confidence between correct and incorrect.

Hypothesis:

$$H_0 : \sigma_1^2 - \sigma_2^2 = 0$$

$$H_1 : \sigma_1^2 - \sigma_2^2 \neq 0$$

Where

σ_1^2 = Variances of the mean of prediction confidence correct.

σ_2^2 = Variances of the mean of prediction confidence incorrect.

Table 35: Result of Test for Equality of Variances of the mean of prediction confidence between correct and incorrect.

	Levene's Test for Equality of Variances	
	F	P-value
Equal variance assumed	4.41	4.01×10^{-2}
* 99.90% confidence intervals (99.90% CI) and p-values from testing ($p \leq 0.001$ was considered statistically significant).		

The test is non-significant, $F = 4.41$, $p = 4.01 \times 10^{-2}$, which indicates that the population variances of correct and incorrect are equal. When equal variances are assumed, the calculation uses pooled variances to use Independent Samples T-Test.

3) Test Statistics

The test statistic for this **Independent Samples T-Test** is denoted t, for equal variances are assumed.

Table 36: Result of Independent Samples T-Test for compare the means of prediction confidence between correct and incorrect group.

Two sample t-test with equal variance				
P - value	t	Mean difference	99.90% Confident Interval of the difference	
			Lower	Upper
1.0×10^{-39}	33.17	39.06	34.98	43.14
*With 99.90% confidence intervals (99.90% CI) and p-values from testing (a one - tailed $p \leq 0.001$ was considered statistically significant).				

B. Compares the means of prediction confidence between correct and incorrect of EfficientNet model

1) Null and Alternative Hypotheses

$$H_0 : \mu_1 = \mu_2$$

$$H_1 : \mu_1 > \mu_2$$

Where

μ_1 = Mean of prediction confidence correct.

μ_2 = Mean of prediction confidence incorrect.

2) The Assumption tests

- There is no relationship of mean of prediction confidence between correct and incorrect.
- Test of Normality: We use **Shapiro-wilk test** to testing normal distribution of mean each prediction confidence.

Prediction confidences correct:

Hypothesis:

H_0 : Mean of prediction confidence correct follows normal distribution.

H_1 : Mean of prediction confidence correct does not follows normal distribution.

Table 37: Result of Test of Normality of the mean of prediction confidence correct.

	Shapiro-wilk	
	W-test statistic	P-value
Correct	0.87	2.0×10^{-3}
* 99.90% confidence intervals (99.90% CI) and p-values from testing ($p \leq 0.001$ was considered statistically significant).		

The test is non-significant, $W = 0.87$, $p = 2.00 \times 10^{-3}$, which indicates that the mean of confidence correct are normally distributed.

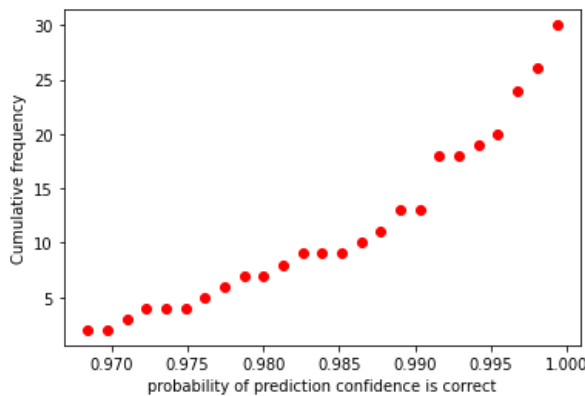


Fig 28: Probability Plots (PP Plot) of the mean prediction confidence correct.

Prediction confidences incorrect:

Hypothesis:

H_0 : Mean of prediction confidence incorrect follows normal distribution.

H_1 : Mean of prediction confidence incorrect does not follows normal distribution.

Table 38: Result of Test of Normality of the mean prediction confidence incorrect.

	Shapiro-wilk	
	W-test statistic	P-value
Incorrect	0.97	0.81
* 99.90% confidence intervals (99.90% CI) and p-values from testing ($p \leq 0.001$ was considered statistically significant).		

The test is non-significant, $W = 0.97$, $p = 0.81$, which indicates that the mean of confidence incorrect are normally distributed.

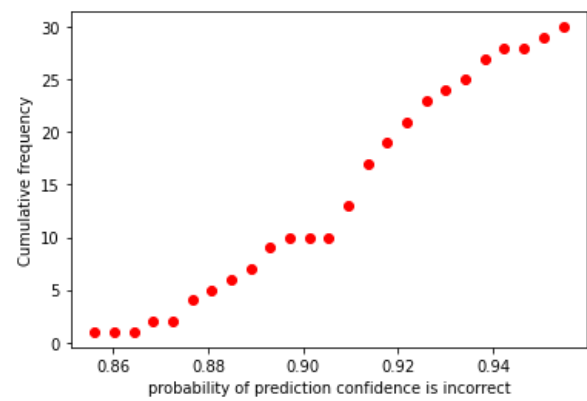


Fig 29: Probability Plots (PP Plot) of the mean prediction confidence incorrect.

- Test of Homogeneity of variances: We use **Levene's Test** to test for the homogeneity of variance of the mean of prediction confidence between correct and incorrect.

Hypothesis

$$H_0 : \sigma_1^2 - \sigma_2^2 = 0$$

$$H_1 : \sigma_1^2 - \sigma_2^2 \neq 0$$

Where

σ_1^2 = Variances of the mean of prediction confidence correct.

σ_2^2 = Variances of the mean of prediction confidence incorrect.

Table 39: Result of Test for Equality of Variances of the mean of prediction confidence between correct and incorrect.

	Levene's Test for Equality of Variances	
	F	P-value
Equal variance not assumed	15.23	2.51×10^{-4}
* 99.90% confidence intervals (99.90% CI) and p-values from testing ($p \leq 0.001$ was considered statistically significant).		

The test is non-significant, $F = 15.23$, $p = 2.51 \times 10^{-4}$, which indicates that the population variances of correct and incorrect are not equal. When equal variances not assumed, the calculation utilizes un-pooled variances to use Independent Samples T-Test.

3) Test Statistics

The test statistic for this **Independent Samples T-Test** is denoted t , for equal variances not assumed.

Table 40: Result of Independent Samples T-Test for compare the means of prediction confidence between correct and incorrect group.

Two sample t-test with unequal variance (Welch's t-test)				
P - value	t	Mean difference	99.90% Confident Interval of the difference	
			Lower	Upper
1.22×10^{-18}	15.74	7.67	5.93	9.41
*With 99.90% confidence intervals (99.90% CI) and p-values from testing (a one - tailed $p \leq 0.001$ was considered statistically significant).				

IV. COMPARES PERFORMANCE OF PARTICIPANTS BETWEEN ASSISTED VS UNASSISTED

We use **Paired Samples T-Test** to compare performance of participants with assisting tool and without assisting tool.

A. Impact of the assisting tool by compare performance of participants in accuracy scores

1) Null and Alternative Hypotheses

$$H_0 : \mu_2 = \mu_1$$

$$H_1 : \mu_2 > \mu_1$$

Where

μ_1 = Mean of accuracy among participants without assisting tool.

μ_2 = Mean of accuracy among participants with assisting tool.

2) The Assumption tests

- There is relationship between accuracy scores among participants with assisting tool and without assisting tool.
- Test of Normality: We use **Shapiro-wilk test** to testing normal distribution of accuracy scores difference between assisted and unassisted.

Hypothesis:

H_0 : Accuracy scores difference between among participants with assisting tool and without the tool follows normal distribution.

H_1 : Accuracy scores difference between among participants with assisting tool and without the tool does not follows normal distribution.

Table 41: Result of Test of Normality of accuracy scores difference between among participants with assisting tool and without the tool.

	Shapiro-wilk	
	W-test statistic	P-value
Assisted - Unassisted	0.90	0.24
* 99.90% confidence intervals (99.90% CI) and p-values from testing ($p \leq 0.001$ was considered statistically significant).		

The test is non-significant, $W = 0.90$, $p = 0.24$, which indicates that the accuracy scores both with assisting tool and without assisting tool are normally distributed.

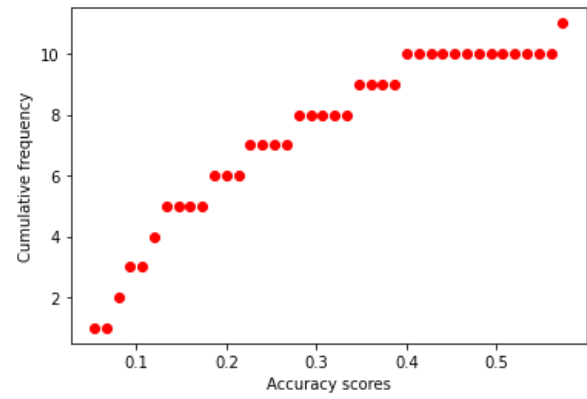


Fig 30: Probability Plots (PP Plot) of accuracy scores difference (assisted - unassisted).

3) Test Statistics

The test statistic for the **Paired Samples T-Test**, denoted t , for compare the means for assisted – unassisted.

Table 42: Result of Paired Samples T-Test between with assisting tool and without assisting tool: accuracy scores.

Paired t-test				
P - value	t	Mean difference	99.90% Confident Interval of the difference	
			Lower	Upper
3.44×10^{-4}	4.83	35.27	1.80	68.75
*With 99.90% confidence intervals (99.90% CI) and p-values from testing (a one - tailed $p \leq 0.001$ was considered statistically significant).				

4) Interval estimates Using T-score with 99.90% CI

Table 43: Result of Interval estimates of accuracy scores using T-score.

Interval estimates using T-score			
Group	Mean of accuracy scores	99.90% Confident Interval	
		Lower	Upper
Assisted	73.52	57.01	90.02
Unassisted	50.00	78.57	21.43

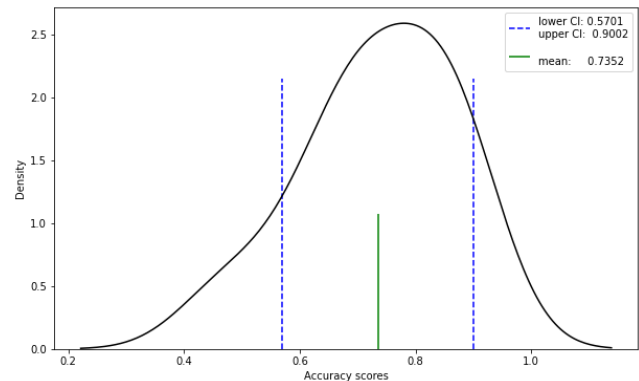


Fig 31: Plot of accuracy scores among participants with assisting tool, t-statistics - Confidence Level = 99.90%.

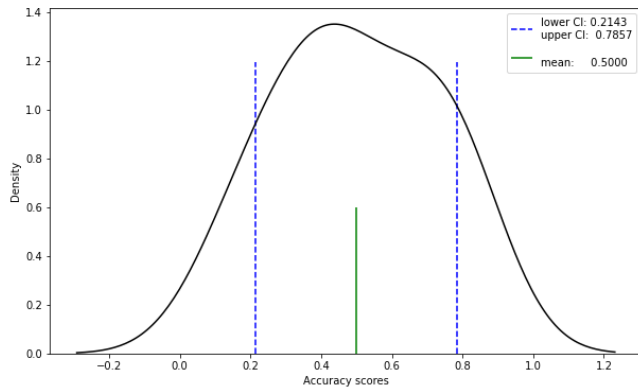


Fig 32: Plot of accuracy scores among participants without assisting tool, t-statistics - Confidence Level = 99.90%.

B. Impact of the assisting tool by compare performance of participants in precision scores

1) Null and Alternative Hypotheses

$$H_0 : \mu_2 = \mu_1$$

$$H_1 : \mu_2 > \mu_1$$

Where

μ_1 = Mean of precision among participants without assisting tool.

μ_2 = Mean of precision among participants with assisting tool.

2) The Assumption tests

- There is relationship between precision scores among participants with assisting tool and without assisting tool.
- Test of Normality: We use **Shapiro-wilk test** to testing normal distribution of precision scores difference between assisted and unassisted.

Hypothesis:

H_0 : Precision scores difference between among participants with assisting tool and without the tool follows normal distribution.

H_1 : Precision scores difference between among participants with assisting tool and without the tool does not follows normal distribution.

Table 44: Result of Test of Normality of precision scores difference between among participants with assisting tool and without the tool.

	Shapiro-wilk	
	W-test statistic	P-value
Assisted - Unassisted	0.95	0.62
* 99.90% confidence intervals (99.90% CI) and p-values from testing ($p \leq 0.001$ was considered statistically significant).		

The test is non-significant, $W = 0.95$, $p = 0.62$, which indicates that the precision scores both with assisting tool and without assisting tool are normally distributed.

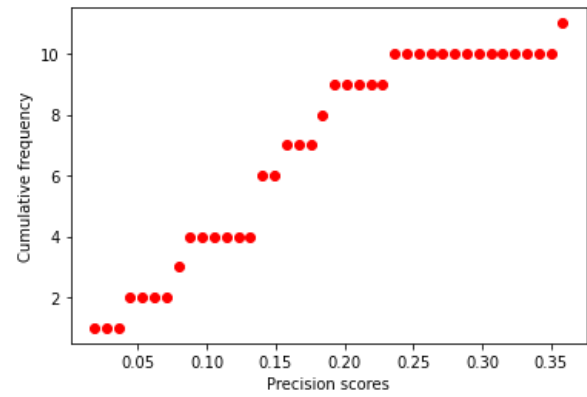


Fig 33: Probability Plots (PP Plot) of precision scores difference (assisted - unassisted).

3) Test Statistics

The test statistic for the **Paired Samples T-Test**, denoted t , for compare the means for assisted – unassisted.

Table 45: Result of Paired Samples T-Test between with assisting tool and without assisting tool: precision scores.

Paired t-test				
P - value	t	Mean difference	99.90% Confident Interval of the difference	
			Lower	Upper
1.58×10^{-4}	5.37	15.39	2.24	28.54
*With 99.90% confidence intervals (99.90% CI) and p-values from testing (a one - tailed $p \leq 0.001$ was considered statistically significant).				

4) Interval estimates Using T-score with 99.90% CI

Table 46: Result of Interval estimates of precision scores using T-score.

Interval estimates using T-score			
Group	Mean of precision scores	99.90% Confident Interval	
		Lower	Upper
Assisted	61.49	42.88	80.10
Unassisted	46.10	25.81	66.38

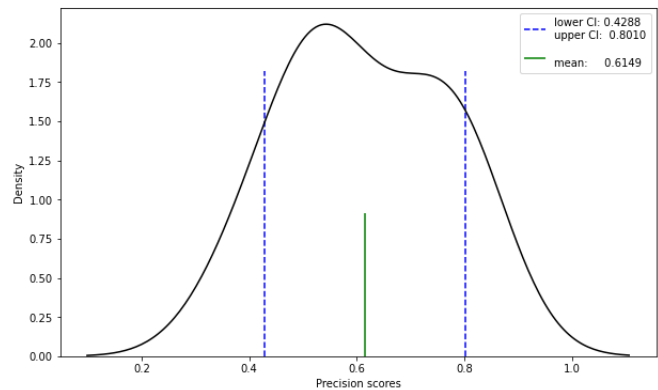


Fig 34: Plot of precision scores among participants with assisting tool, t-statistics - Confidence Level = 99.90%.

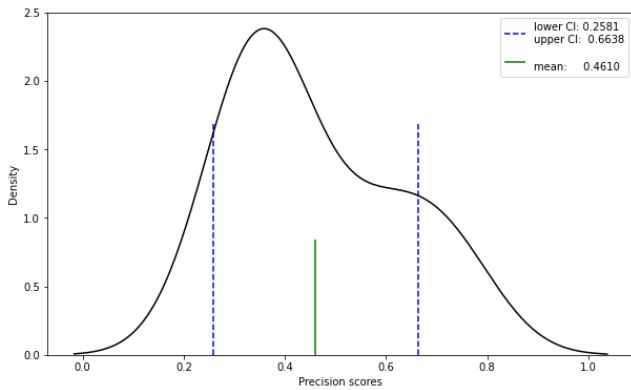


Fig 35: Plot of precision scores among participants without assisting tool, t-statistics - Confidence Level = 99.90%.

C. Impact of the assisting tool by compare performance of participants in recall scores

1) Null and Alternative Hypotheses

$$H_0 : \mu_2 = \mu_1$$

$$H_1 : \mu_2 > \mu_1$$

Where

μ_1 = Mean of recall among participants without assisting tool.

μ_2 = Mean of recall among participants with assisting tool.

2) The Assumption tests

- There is relationship between recall scores among participants with assisting tool and without assisting tool.
- Test of Normality: We use **Shapiro-wilk test** to testing normal distribution of recall scores difference between assisted and unassisted.

Hypothesis:

H_0 : Recall scores difference between among participants with assisting tool and without the tool follows normal distribution.

H_1 : Recall scores difference between among participants with assisting tool and without the tool does not follows normal distribution.

Table 47: Result of Test of Normality of recall scores difference between among participants with assisting tool and without the tool.

	Shapiro-wilk	
	W-test statistic	P-value
Assisted - Unassisted	0.94	0.57
* 99.90% confidence intervals (99.90% CI) and p-values from testing ($p \leq 0.001$ was considered statistically significant).		

The test is non-significant, $W = 0.94$, $p = 0.57$, which indicates that the recall scores both with assisting tool and without assisting tool are normally distributed.

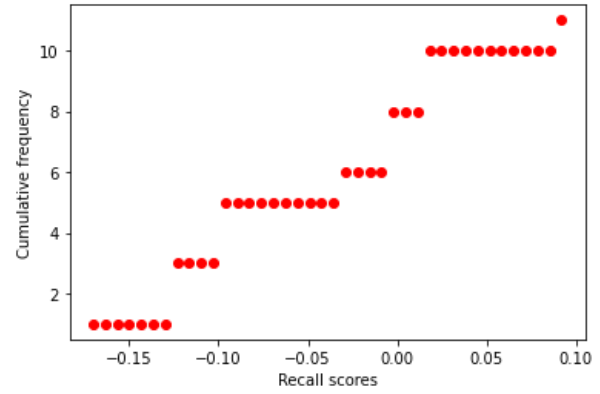


Fig 36: Probability Plots (PP Plot) of recall scores difference (assisted - unassisted).

3) Test Statistics

The test statistic for the **Paired Samples T-Test**, denoted t , for compare the means for assisted – unassisted.

Table 48: Result of Paired Samples T-Test between with assisting tool and without assisting tool: recall scores.

Paired t-test				
P - value	t	Mean difference	99.90% Confident Interval of the difference	
			Lower	Upper
0.05	-1.79	-4.33	-15.42	6.77
*With 99.90% confidence intervals (99.90% CI) and p-values from testing (a one - tailed $p \leq 0.001$ was considered statistically significant).				

4) Interval estimates Using T-score with 99.90% CI

Table 49: Result of Interval estimates of recall scores using T-score.

Interval estimates using T-score			
Group	Mean of recall scores	99.90% Confident Interval	
		Lower	Upper
Assisted	88.31	79.34	97.28
Unassisted	92.64	85.30	99.98

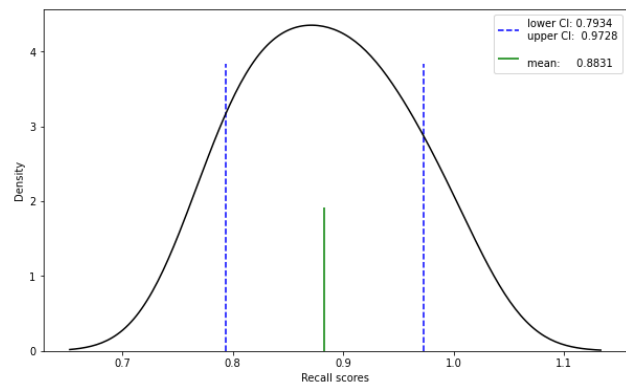


Fig 37: Plot of recall scores among participants with assisting tool, t-statistics - Confidence Level = 99.90%.

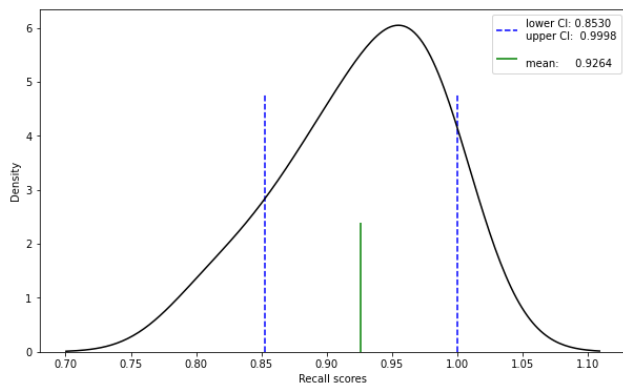


Fig 38: Plot of recall scores among participants without assisting tool, t-statistics - Confidence Level = 99.90%.

The test is non-significant, $W = 0.94$, $p = 0.55$, which indicates that the accuracy scores difference between the first round of experiment and the second round of experiment follows normal distribution.

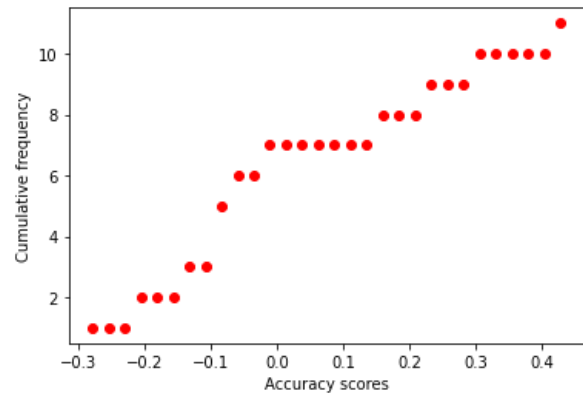


Fig 39: Probability Plots (PP Plot) of accuracy scores difference (second experiment – first experiment).

V. THE PERFORMANCE OF THE PARTICIPANTS BETWEEN THE FIRST ROUND OF EXPERIMENT AND THE SECOND ROUND OF EXPERIMENT

We use **Paired Samples T-Test** to compare accuracy between the first round of experiment and the second round of experiment of the participants.

5.1 Null and Alternative Hypotheses

$$H_0 : \mu_2 - \mu_1 = 0$$

$$H_1 : \mu_2 - \mu_1 \neq 0$$

Where

μ_1 = Mean of accuracy first round of experiment.

μ_2 = Mean of accuracy second round of experiment.

5.2 The Assumption tests

1) There is relationship of accuracy scores the rounds of experiments, between the first session and the second session.

2) Test of Normality: We use **Shapiro-wilk test** to testing normal distribution between the Accuracy scores of 11 participants on the first and the second sessions.

Hypothesis:

H_0 : Accuracy scores difference between the first round of experiment and the second round of experiment follows normal distribution.

H_1 : Accuracy scores difference between the first round of experiment and the second round of experiment does not follows normal distribution.

Table 50: Result of Test of Normality of accuracy scores difference between of participants between the first round of experiment and the second round.

	Shapiro-wilk	
	W-test statistic	P-value
Second experiment – First experiment	0.94	0.55
* 99.90% confidence intervals (99.90% CI) and p-values from testing ($p \leq 0.001$ was considered statistically significant).		

5.3 Test Statistics

The test statistic for the **Paired Samples T-Test**, denoted t , for compare the means of accuracy for the first and the second sessions.

Table 51: Result of Paired Samples T-Test for compare the means of accuracy the first round of experiment and the second round of experiment.

Paired t-test				
P - value	t	Mean difference	99.90% Confident Interval of the difference	
			Lower	Upper
0.57	0.59	4.00	27.04	35.04
*With 99.90% confidence intervals (99.90% CI) and p-values from testing (a two-tailed $p \leq 0.001$ was considered statistically significant).				

5.4 Interval estimates Using T-score with 99.90% CI

Table 52: Result of Interval estimates of accuracy scores using T-score.

Interval estimates using T-score			
Group	Mean of accuracy scores	99.90% Confident Interval	
		Lower	Upper
First experiment	68.24	38.14	98.34
Second experiment	72.24	47.52	96.97

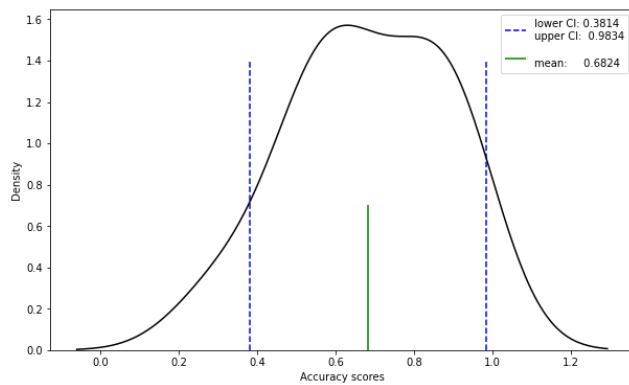


Fig 40: Plot of accuracy scores of participants on the first experiment, t-statistics - Confidence Level = 99.90%.

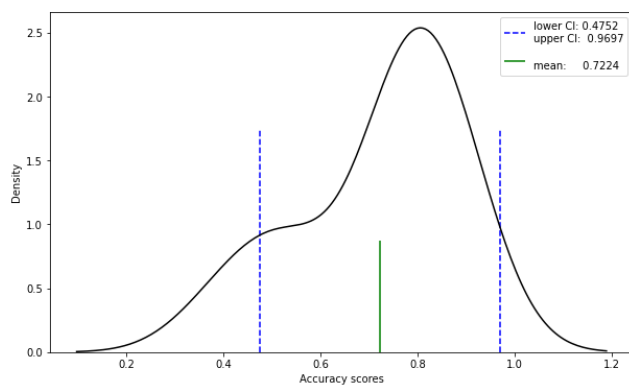


Fig 41: Plot of accuracy scores of participants on the second experiment, t-statistics - Confidence Level = 99.90%.

VI. INFLUENCE OF AI SUGGESTION ON PARTICIPANT DECISIONS WHEN ASSISTED/UNASSISTED

We use **Paired Samples T-Test** to compare similarity scores between AI suggestion (prediction) and the final decision of the participants when assisted/unassisted.

6.1 Null and Alternative Hypotheses

$$H_0 : \mu_2 = \mu_1$$

$$H_1 : \mu_2 > \mu_1$$

Where

μ_1 = Mean of similarity between AI suggestion and participant decisions without assisting tool.

μ_2 = Mean of similarity between AI suggestion and participant decisions with assisting tool.

6.2 The Assumption tests

1) There is relationship of similarity scores between AI suggestion and decision of 11 participants when assisted/unassisted.

2) Test of Normality: We use **Shapiro-wilk test** to testing normal distribution between the similarity scores between AI suggestion and participant decisions when assisted/unassisted.

Hypothesis

H_0 : Similarity scores difference between AI suggestion and participant decisions when assisted/unassisted follows normal distribution.

H_1 : Similarity scores difference between AI suggestion and participant decisions when assisted/unassisted does not follows normal distribution.

Table 53: Result of Test of Normality of similarity scores difference between AI suggestion and participant decisions when assisted/unassisted.

	Shapiro-wilk	
	W-test statistic	P-value
Assisted - Unassisted	0.94	0.49
* 99.90% confidence intervals (99.90% CI) and p-values from testing ($p \leq 0.001$ was considered statistically significant).		

The test is non-significant, $W = 0.94$, $p = 0.49$, which indicates that the similarity scores difference between AI suggestion and participant decisions when assisted/unassisted follows normal distribution.

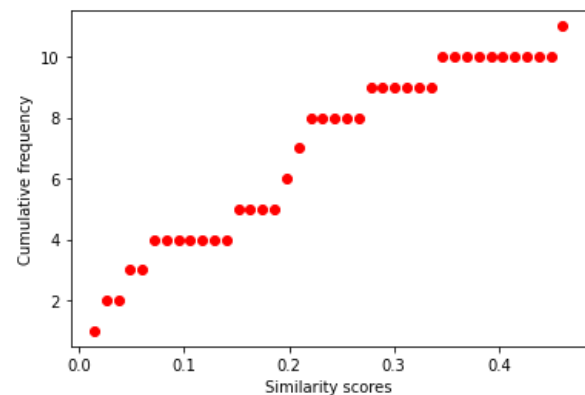


Fig 42: Probability Plots (PP Plot) of similarity scores difference between AI suggestion and participant decisions (assisted - unassisted).

6.3 Test Statistics

The test statistic for the **Paired Samples T-Test**, denoted t , for compare the means of similarity for participant decisions when assisted/unassisted.

Table 54: Result of Paired Samples T-Test for compare the means of similarity between AI suggestion and participant decisions when assisted/unassisted.

Paired t-test				
P - value	t	Mean difference	99.90% Confident Interval of the difference	
			Lower	Upper
6.90×10^{-4}	4.38	18.78	-0.89	38.47
*With 99.90% confidence intervals (99.90% CI) and p-values from testing (a one - tailed $p \leq 0.001$ was considered statistically significant).				

6.4 Interval estimates Using T-score with 99.90% CI

Table 55: Result of Interval estimates of similarity scores using T-score.

Interval estimates using T-score			
Group	Mean of similarity scores	99.90% Confident Interval	
		Lower	Upper
Assisted	77.64	63.47	91.81
Unassisted	58.85	34.07	83.63

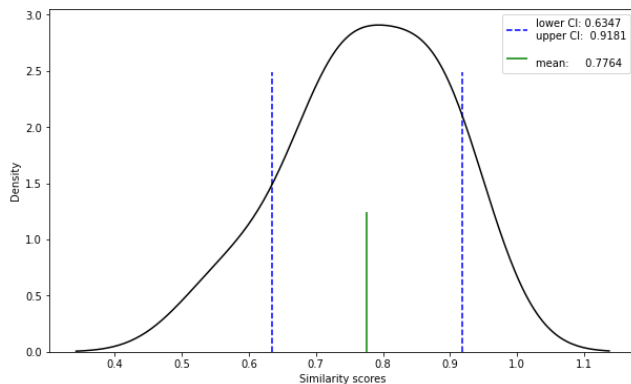


Fig 43: Plot of similarity scores between AI suggestion and participant decisions when assisted, t-statistics - Confidence Level = 99.90%.

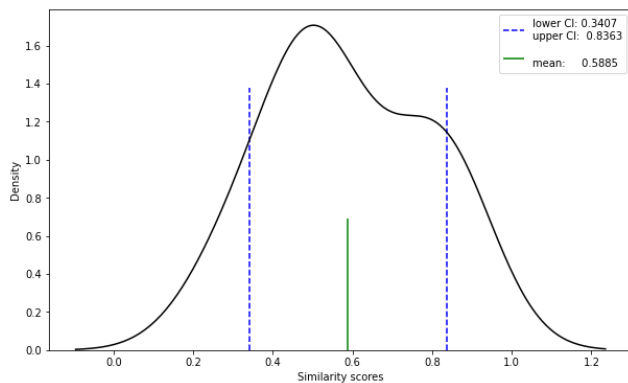


Fig 44: Plot of similarity scores between AI suggestion and participant decisions when unassisted, t-statistics - Confidence Level = 99.90%.

VII. COMPARE THE RELATIONSHIP BETWEEN HIGH-LOW PREDICTION CONFIDENCE AND SIMILARITY OF THE PARTICIPANT ANSWER

We use **Pearson Chi-Square test** to hypothesis testing correlation between high-low prediction confidence (confidence ≤ 50 and confidence > 50) and similarity of the participant answer to the prediction suggestion suggested.

7.1 Our cross-tabulation table

Table 56: Cross tabulation between high-low prediction confidence and similarity of the participant answer to the

prediction suggested.

Prediction confidence	The answer of participant		Total
	Does not have similar answer	Have similar answer	
High	331	956	1,287
Low	181	182	363
Total	512	1,138	1,650

7.2 Null and Alternative Hypotheses

H_0 : Prediction confidence is not associated with the answer of participant.

H_1 : Prediction confidence is associated with the answer of participant.

7.3 The Assumption tests

- 1) Prediction confidence and the answer of participant were collected independently of each other.
- 2) Whole expected cell counts greater than 10.
We can be checked by looking at the expected frequency table.

Table 57: Expected frequency table between high-low prediction confidence and similarity of the participant answer to the prediction suggested.

Prediction confidence	The answer of participant	
	Does not have similar answer	Have similar answer
High	399.36	887.64
Low	112.64	250.36

7.4 Test Statistics

The test statistic for the **Chi-Square Test of Independence** is denoted χ^2 , the research question is the following, is there a relationship between prediction confidence and the answer of participant.

Table 58: Result of Chi-Square Test of Independence between prediction confidence and the answer of participant.

	Value	P - value
Pearson Chi-Square	76.00	2.84×10^{-18}
*With 99.90% confidence intervals (99.90% CI) and p-values from testing (a two - tailed $p \leq 0.001$ was considered statistically significant).		

VIII. COMPARE THE RELATIONSHIP BETWEEN CORRECT-INCORRECT ROI AND THE PARTICIPANT DECISIONS.

We use **Pearson Chi-Square test** to hypothesis testing correlation between the decisions when IoU of the GradCam and the ROI are greater than 0.8 (correct) and the decisions when IoU of the GradCam and the ROI are less than 0.3 (incorrect).

8.1 Our cross-tabulation table

Table 59: Cross tabulation between correct - incorrect IoU and decisions of the participant.

Rating of IOU	The decisions of participant		Total
	Does not have similar decisions	Have similar decisions	
Correct	2	20	22
Incorrect	96	69	165
Total	98	89	187

8.2 Null and Alternative Hypotheses

H_0 : IoU of the GradCam and the ROI is not associated with the decisions of the participant.

H_1 : IoU of the GradCam and the ROI is associated with the decisions of the participant.

8.3 The Assumption tests

- 1) IoU value and decisions of the participant were collected independently of each other.
- 2) Whole expected cell counts greater than 10.
We can be checked by looking at the expected frequency table.

Table 60: Expected frequency table between correct - incorrect IoU and decisions of the participant.

Rating of IOU	The decisions of participant	
	Does not have similar decisions	Have similar decisions
Correct	11.53	10.47
Incorrect	86.47	78.53

8.4 Test Statistics

The test statistic for the **Chi-Square Test of Independence** is denoted χ^2 , the research question is the following, is there a relationship between IoU of the GradCam and the ROI and the decisions of the participant.

Table 61: Result of Chi-Square Test of Independence between correct - incorrect IoU and decisions of the participant.

	Value	P - value
Pearson Chi-Square	16.84	4.07×10^{-5}
*With 99.90% confidence intervals (99.90% CI) and p-values from testing (a two - tailed $p \leq 0.001$ was considered statistically significant).		

IX. COMPARE THE RELATIONSHIP BETWEEN CORRECT-INCORRECT VIEWING ANGLE PREDICTION AND THE PARTICIPANT DECISIONS.

We use **Pearson Chi-Square test** to hypothesis testing correlation between the decisions when the viewing angle predictions are correct and the decisions when the viewing angle predictions are incorrect.

9.1 Our cross-tabulation table

Table 62: Cross tabulation between correct – incorrect viewing angle predictions and decisions of the participant.

Viewing angle predictions	The decisions of participant		Total
	Does not have similar decisions	Have similar decisions	
Correct	299	779	1,078

Viewing angle predictions	The decisions of participant		Total
	Does not have similar decisions	Have similar decisions	
Incorrect	196	376	572
Total	495	1,155	1,650

9.2 Null and Alternative Hypotheses

H_0 : Viewing angle predictions is not associated with the decisions of the participant.

H_1 : Viewing angle predictions is associated with the decisions of the participant.

9.3 The Assumption tests

- 1) Viewing angle predictions and decisions of the participant were collected independently of each other.
- 2) Whole expected cell counts greater than 10.
We can be checked by looking at the expected frequency table.

Table 63: Expected frequency table between correct - incorrect viewing angle predictions and decisions of the participant.

Viewing angle predictions	The decisions of participant	
	Does not have similar decisions	Have similar decisions
Correct	323.40	754.60
Incorrect	171.60	400.40

9.4 Test Statistics

The test statistic for the **Chi-Square Test of Independence** is denoted χ^2 , the research question is the following, is there a relationship between viewing angle predictions and the decisions of the participant.

Table 64: Result of Chi-Square Test of Independence between correct - incorrect viewing angle predictions and decisions of the participant.

	Value	P - value
Pearson Chi-Square	7.28	7.00×10^{-3}
*With 99.90% confidence intervals (99.90% CI) and p-values from testing (a two - tailed $p \leq 0.001$ was considered statistically significant).		

X. INFLUENCE OF TOP-3 PREDICTION ON PARTICIPANT DECISIONS

We use **Paired Samples T-Test** to compare similarity scores between the participant decisions versus the model top second predictions or the model top third predictions, assisted and unassisted.

10.1 Null and Alternative Hypotheses

$H_0 : \mu_2 = \mu_1$

$H_1 : \mu_2 > \mu_1$

Where

μ_1 = Mean of similarity between top-3 prediction and participant decisions without assisting tool.

μ_2 = Mean of similarity between top-3 prediction and participant decisions with assisting tool.

10.2 The Assumption tests

- 1) There is relationship of similarity scores between top-3 prediction and decision of 11 participants when assisted/unassisted.
- 2) Test of Normality: We use **Shapiro-wilk test** to testing normal distribution between the similarity scores between top-3 prediction and participant decisions when assisted/unassisted.

Hypothesis:

H_0 : Similarity scores difference between top-3 prediction and participant decisions when assisted/unassisted follows normal distribution.

H_1 : Similarity scores difference between top-3 prediction and participant decisions when assisted/unassisted does not follows normal distribution.

Table 65: Result of Test of Normality of similarity scores difference between top-3 prediction and participant decisions when assisted/unassisted.

	Shapiro-wilk	
	W-test statistic	P-value
Assisted - Unassisted	0.92	0.31
* 99.90% confidence intervals (99.90% CI) and p-values from testing ($p \leq 0.001$ was considered statistically significant).		

The test is non-significant, $W = 0.92$, $p = 0.31$, which indicates that the similarity scores difference between top-3 prediction and participant decisions when assisted/unassisted follows normal distribution.

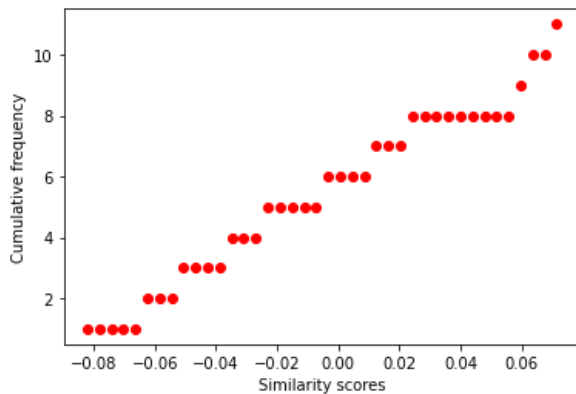


Fig 45: Probability Plots (PP Plot) of similarity scores difference between top-3 prediction and participant decisions (assisted - unassisted).

10.3 Test Statistics

The test statistic for the **Paired Samples T-Test**, denoted t , for compare the means of similarity for participant decisions when assisted/unassisted.

Table 66: Result of Paired Samples T-Test for compare the means of similarity between top-3 prediction and participant decisions when assisted/unassisted.

Paired t-test				
P - value	t	Mean difference	99.90% Confident Interval of the difference	
			Lower	Upper
0.50	0.00	-2.52×10^{-16}	-8.61	8.61
*With 99.90% confidence intervals (99.90% CI) and p-values from testing (a one - tailed $p \leq 0.001$ was considered statistically significant).				

10.4 Interval estimates Using T-score with 99.90% CI

Table 67: Result of Interval estimates of similarity scores using T-score.

Interval estimates using T-score			
Group	Mean of similarity score	99.90% Confident Interval	
		Lower	Upper
Assisted	14.73	8.17	21.28
Unassisted	14.73	10.33	19.13

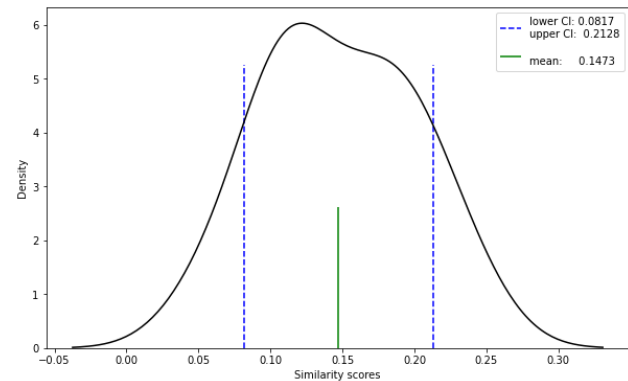


Fig 46: Plot of similarity scores between top-3 prediction and participant decisions when assisted, t-statistics - Confidence Level = 99.90%.

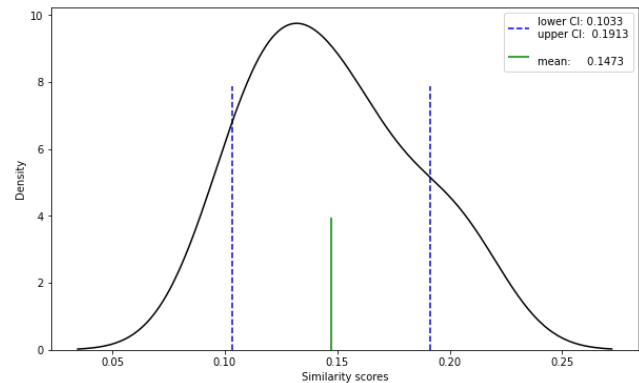


Fig 47: Plot of similarity scores between top-3 prediction and participant decisions when unassisted, t-statistics - Confidence Level = 99.90%.

XI. CONFUSION MATRICES OF THE PERFORMANCE OF PARTICIPANTS ON DIFFERENT ABNORMALITIES.

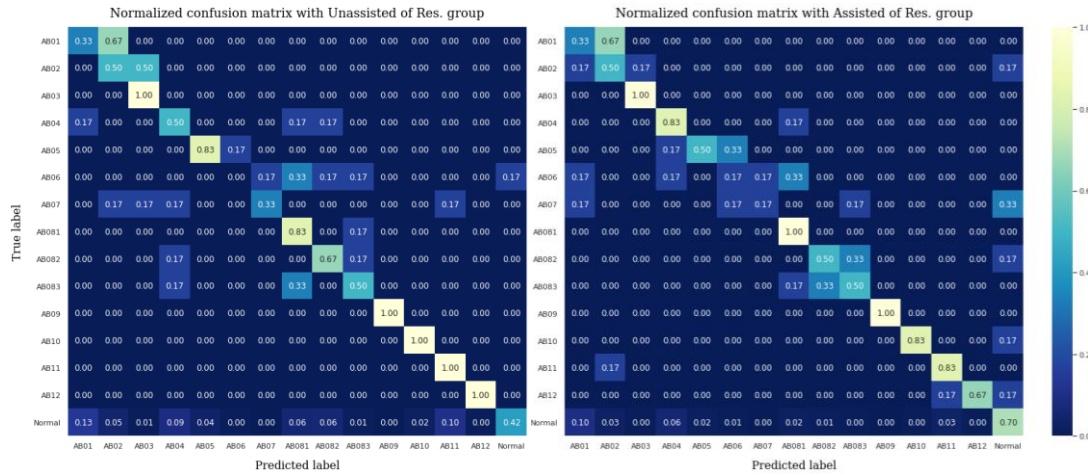


Fig 48: The confusion matrix of the performance of the radiologist residence group without the assisting tool (left) and with assisting tool (right), the numbers are row-wise normalization.

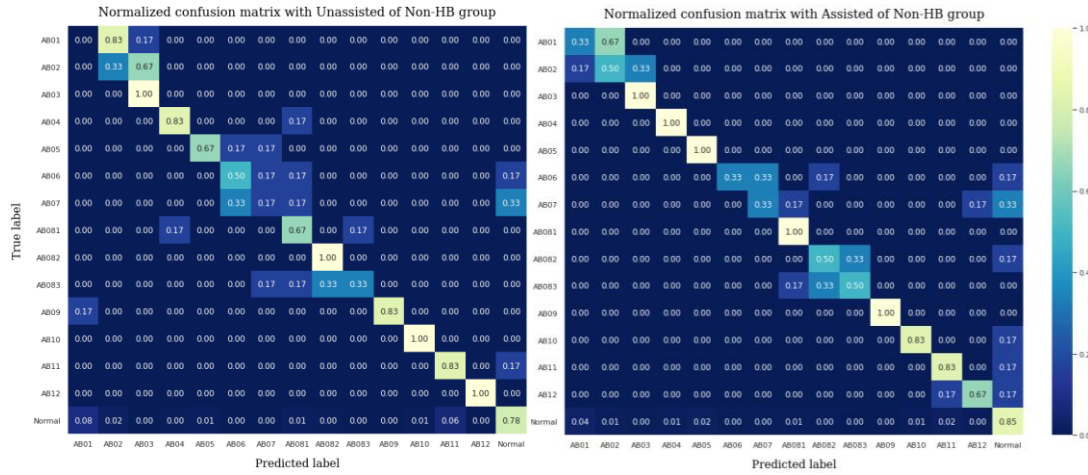


Fig 49: The confusion matrix of the performance of the non-hepatobiliary radiologist group without the assisting tool (left) and with assisting tool (right), the numbers are row-wise normalization.

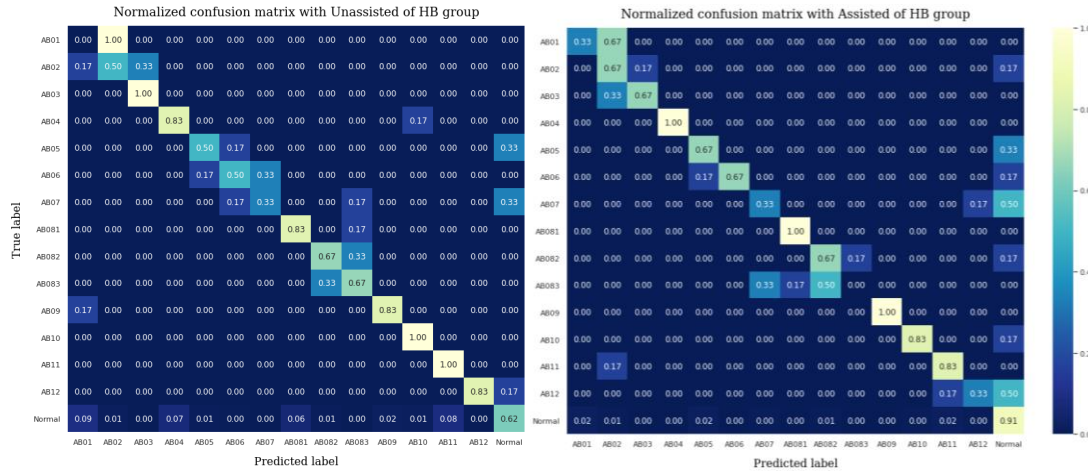


Fig 50: The confusion matrix of the performance of the hepatobiliary radiologist group without the assisting tool (left) and with assisting tool (right), the numbers are row-wise normalization.