



**CALIFORNIA STATE SCIENCE FAIR  
2014 PROJECT SUMMARY**

<b>Name(s)</b> Kevin A. Hsieh; Abraham N. Razzak	<b>Project Number</b>  34955
<b>Project Title</b> <b>Optimizing Quantization Matrix Scale Factor and Pixel Density for Effective Optical Character Recognition in JPEG Images</b>	
<b>Objectives/Goals</b> The purpose of this experiment was to determine the effect that changes in quantization matrix scale factor (abbreviated "S-factor") and pixel density have on the effectiveness of optical character recognition (OCR) in JPEG images, and to create and use a script for finding optimal values for those variables under a given set of conditions. It was hypothesized that OCR effectiveness increases (1) as pixel density increases and (2) as S-factor decreases, because increasing pixel density improves clarity, but increasing S-factor introduces compression artifacts. <b>Abstract</b> <b>Methods/Materials</b> This experiment requires a computer capable of running batch files. Other software programs required included IrfanView and Tesseract OCR. In order to test the effect that the manipulated variables of S-factor and pixel density have on OCR, a script based on batch, a scripting language used to coordinate programs with large numbers of files, was created. The script was run three times, on three different sets of words, for three trials. This resulted in 29,700 data points, which were then tabulated and converted from correctly transcribed word count to a percentage representing OCR effectiveness. Visuals were created from the data for analysis purposes. <b>Results</b> In most trials, OCR recognized at least some words when pixel density was at least 40 ppcm, and recognized all of the words once the pixel density reached approximately 75 ppcm. A large amount of deviation was observed between 40 and 75 ppcm, but overall, OCR effectiveness improved as pixel density increased. The effect of Q-factor was unclear at first; data points had to be re-expressed in terms of S-factor for statistical analysis. It was observed that S-factor and OCR effectiveness were inversely proportional. <b>Conclusions/Discussion</b> The results supported both parts of the hypothesis and the associated justifications. This experiment assists in understanding what combination of S-values/Q-values and pixel density produce useful OCR results, a concept that is relevant when determining optimal compression ratios for images intended to be processed with OCR, a process common in modern smartphones. Under the conditions of this experiment, a JPEG image at about 75 ppcm with an S-value of about 1.67 (Q = 30) would have allowed for maximum OCR effectiveness at a small file size.	
<b>Summary Statement</b> This project determined the effect that changes in S-factor and pixel density have on the effectiveness of OCR in JPEG images, and created and used a script for finding optimal values for those variables under a given set of conditions.	
<b>Help Received</b> Mr. Antrim advised and oversaw progress; Mr. Tipper helped with statistical analysis; Parents provided equipment and locations for testing	