

Leveraging Cognitive Factors in Securing WWW with CAPTCHA

Amalia Rusu and Rebecca Docimo
Department of Software Engineering



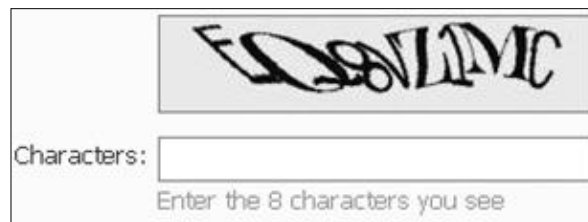
Adrian Rusu
Department of Computer Science



Presented by: Amalia Rusu

Introduction / Motivation

- CAPTCHAs (*Completely Automated Public Turing Tests to Tell Computers and Humans Apart*) widespread on WWW to prevent illegitimate access to web resources by programs (aka “bots” or machines). Best current solution, though important insights can also be gained if broken. CAPTCHAs are often made harder against machines at the expense of human solvability.
- Paper presents the development of secure but *usable* new system to protect Web applications.
 - Tree-Based Handwritten CAPTCHA protects Web applications while ensuring human legibility by utilizing principles of human cognition (Gestalt principles and Geon theory) to motivate transformations. or precise changes, to handwritten words. Transformations added to already difficult to segment handwriting presents challenges to machines. Handwriting interpretation, despite condition, is a trivial task for humans.
 - Transformed handwritten words are added to a tree structure to make a multi-layer drawing easy for humans to interpret but hard for machines due to recognizer weaknesses.
 - Insights are offered into fields with open problems such as Handwriting Recognition, Image Analysis, Human Computer Interaction, etc.
 - Motivated by weaknesses of other CAPTCHAs, open problems in machine recognition of handwriting and complex drawings, and previous success with handwriting only CAPTCHA.



A typical CAPTCHA on the WWW

Leveraging Cognitive Factors in Securing WWW with CAPTCHA

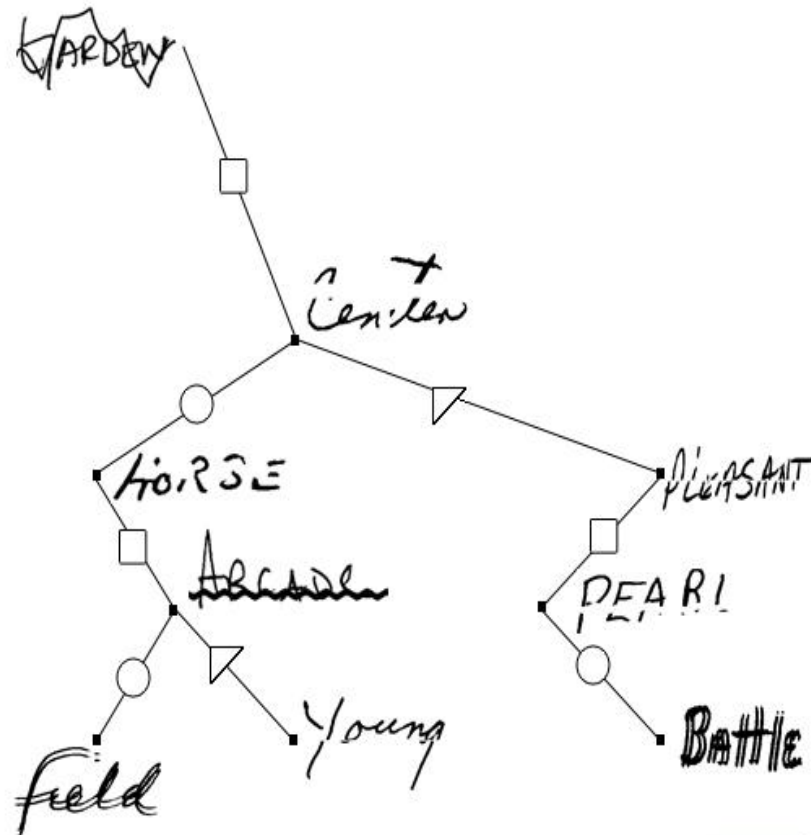
Tree-Based Handwritten CAPTCHA

Binary tree drawing:

- rendered on the fly, variable # edges and nodes
- randomly drawn and transformed handwritten words are context free and placed at each node
- drawn symbols may be added along edges

Solving the challenge:

- all words must be segmented out and interpreted to answer
- correct word must be typed in exactly to pass



Which word is connected to center by a line marked with a circle?

I Can't See the Tree!

Submit

Leveraging Cognitive Factors in Securing WWW with CAPTCHA

Amalia Rusu and Rebecca Docimo, Fairfield University, Fairfield, CT and Adrian Rusu, Rowan University, Glassboro, NJ

Technical Approach

Key elements of Tree-Based Handwritten CAPTCHA:

1. Handwriting Recognition

- Key feature of CAPTCHA is handwriting. Machine skill at segmenting open-context handwriting open problem, while human skill is excellent starting from young age irrespective of writing variability or quality due to cognitive strengths.
- Transformations applied to handwriting further obscure machine segmentation while Gestalt principles and Geon theory preserve human readability.

2. Graphics Recognition

- Trees are simpler structures than graphs, and basic structures for information visualization. Algorithms to visualize trees of any size in a clear manner have been developed.
- Tree structure leverages open problems in Document Image Analysis/Graphics Recognition, for example, lack of cross-domain recognition tools to interpret general drawings without specific primitive elements.
- In all cases, both interpreting the tree and transformed handwriting samples need to be recognized to solve CAPTCHA.

Technical Approach

Key elements of Tree-Based Handwritten CAPTCHA:

4. Geon Theory of Pattern Recognition to Transform Handwriting

- Describes core components which must be left in visual object to remain understandable to a human.
- Geons are simple structures such as cylinders, wedges, cubes, spheres. Edges and intersections are components critical for understanding.
- In general, object recognition is easy for humans if geons can be visualized and recognized.



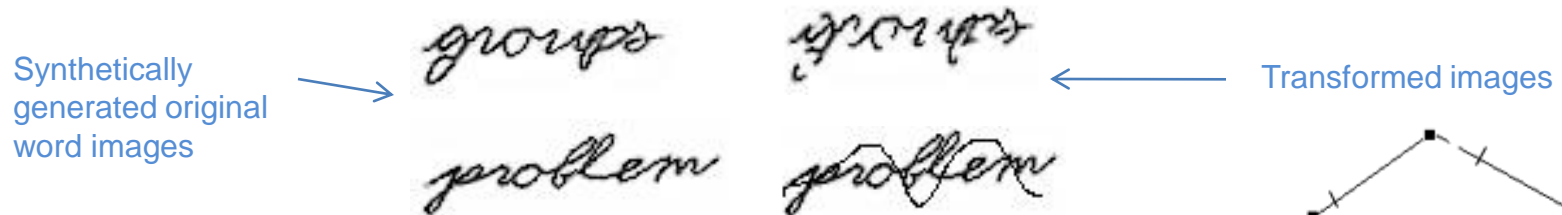
Object with Geons Intact



Object with Missing Geons

Specific Transformations

Principles translated into actual transformations applied to handwritten words. Tree itself or tree elements could likewise be transformed.



- a) add occlusions by waves (background color). **Principles:** closure, proximity, continuity.
- b) add circular, linear or rectangular occlusions (background color). **Principles:** closure, proximity, continuity, familiarity.
- c) use empty/broken objects, rough contours or fragmentation. **Principles:** closure, proximity, continuity and figure-ground.
- d) add occlusions (foreground color). **Principles:** familiarity and figure-ground.
- e) split image into parts and displace or spread as mosaic effect. **Principles:** closure, proximity, continuity, symmetry.
- f) create horizontal/vertical overlaps. **Principles:** proximity, symmetry, familiarity, continuity, figure-ground.
- g) add extra strokes. **Principles:** familiarity, figure-ground.
- h) change orientation, stretch/compress. **Principles:** memory, internal metrics, familiarity of objects.



Leveraging Cognitive Factors in Securing WWW with CAPTCHA

Generating Tree-Based Handwritten CAPTCHA

- Tree generation algorithm uses randomness to automatically draw trees overlaid with synthetically generated and transformed handwritten images as follows:
 1. Random number nodes created.
 2. Once nodes determined, binary tree structure built.
 3. Randomly selected scaling and sizing applied to ensure drawing variability.
 4. Randomly generated synthetic handwritten text placed next to nodes.
 5. Optionally, for each tree branch, randomly drawn symbol placed in middle.
- To complete generation, program selects one or more nodes, or edges, about which to ask a randomly chosen question and passes handwritten images and truth words to the verifier.
- To solve challenge, the tree and handwriting samples must be segmented and interpreted. If correct answer is given by user, challenge is passed and access to Web resource is given, otherwise user fails and new challenge is rendered.

HIP System Evaluation

- Experimental tests conducted with human subjects and on state-of-the-art handwriting recognizers Word Model Recognizer (WMR), Character Model Recognizer (CMR), and Accuscript (HMM).
- To create a “fair” test for machines, word recognizers were assisted with lexicons that contained all the truth words of the test images.
- Scanned handwritten image samples of US city names (lexicon approximately 40,000), readily available from postal applications were used to compare against synthetically generated images.
- Actual CAPTCHA application has no context-specific dictionary so lexicon would be far larger, with machine recognition drastically reduced.

HIP System Evaluation

Machine Testing

- Tests run using human-written scanned samples or synthetically generated samples.
- Transformations based on Gestalt principles and Geon theory, as described earlier, were applied. Parameter values for transformations carefully selected and successively applied.
- For scanned samples, machine recognition rates were very low (1.21% - 5.74%) even with relatively small lexicon of 40,000. Accuracies approached 0% for particular transformations such as word overlaps, adding extra strokes, or fragmentations, conversely the most legible for humans.
- For 300 synthetic samples, recognition was only 0.3% - 1.00% even with small lexicon.
- No commercial product we are aware of can fully interact with the Tree-Based Handwritten CAPTCHA. For further testing, custom scripts may be created.

HIP System Evaluation

Usability Testing

- Focus- to understand viability of CAPTCHA from user experience and solvability perspective. Compared results for human-written samples to synthetically generated images, and handwriting-only CAPTCHA to Tree-Based Handwritten CAPTCHA.
- **Handwriting only CAPTCHA:**
 - 10 handwritten words for each of 9 transformation types considered given to 9 volunteers.
 - Images chosen from scanned US city name samples and lexicon provided to recognizers to ensure fair comparison to machine testing.
 - Overall success - 80%, most errors coming from poor quality scans. Synthetic samples tested with recognition of 80% or better.
- **Tree-Based Handwritten CAPTCHA:**
 - 15 volunteers, 30% non-native English speakers. 190 challenges completed and rated from 1-5 (1 least difficult, 5 most difficult).
 - 80.6% success rate, no less than handwriting-only CAPTCHA.
 - Most common rating was 2. Even trials rated 4-5 often correctly solved, suggests cognitive principles helped users “fill in the blanks”.

Conclusions

- Early experimental tests show that Tree-Based Handwritten CAPTCHA was readily solvable by humans but not by machines and is thus an effective solution to protect Web resources.
- Our CAPTCHA also provides fertile ground for work on important problems in AI, Image Analysis and others.
- Next steps:
 - Additional testing will be run using larger lexicon of words and all synthetic samples and a bigger, more diverse test group. We expect machine recognition to further decrease while human legibility increases.
 - Additional metrics will be collected on time to solve as well as on machine processing load.
 - Custom scripts may be created to attempt to interact with the CAPTCHA as a whole.
 - Extensions to our approach may include combining handwritten text images with images of objects and exploring accessibility options. Transformations may be made to the tree as well.