# Scale and Translation Invariance for Novel Objects in Human Vision

**Vijeta Kumari**

Department of Computer Science, Goel Institute of Technology and Management, Lucknow, India
E-mail: vijetakumari4778@gmail.com

**Abstract:** Understanding invariance in human visual recognition is challenging. Our study investigated the tolerance to scale and position changes during one-shot learning. Non-Japanese individuals showed impressive scale invariance but limited position invariance when recognising Japanese letters. These findings contribute to our understanding of the complex mechanisms underlying human visual recognition. Our study employed both experimental data analysis and computational modelling to investigate the neural computations involved in invariant object recognition. The findings strongly emphasise the necessity of integrating inherent scale-invariance mechanisms into neural network models. This entails incorporating distinct scale channels and accounting for the diverse receptive field sizes and sampling densities of neurons as a function of eccentricity. By adopting this approach, we gain valuable insights into the complex neural processes underlying invariance in human visual recognition, driving future advancements in the field. Through a combination of psychophysical experiments and extensive simulations, our study presents compelling evidence that highlights the striking disparity between the computational strategies employed by the human visual system and prevailing deep learning architectures. Notably, we shed light on the remarkable data efficiency innate to human vision, enabling the extraordinary achievement of successful one-shot learning—a feat that continues to elude conventional deep learning models. Additionally, our investigation emphasises the heightened importance of eye movements as a vital mechanism for object recognition within the human visual system. These distinct characteristics unravel fresh insights into the intricate interplay between human cognition and visual perception, fostering a deeper understanding of the inherent complexities involved in visual processing and paving the way for innovative avenues of inquiry. Our study enhances our understanding of scale and position invariance in human vision, shedding light on the underlying neural mechanisms. The implications emphasise the importance of integrating scale-invariant mechanisms into neural network models and recognising the critical role of eye movements in achieving object recognition invariance. These insights drive further investigations and push the boundaries of knowledge in visual perception.

**Keywords:** Deep Learning, ENN, Learning Algorithm, CNN, Machine Learning.

**Introduction-** The study of scale and position invariance in human vision has significant implications for computer vision, AI, and cognitive neuroscience. It enhances object recognition algorithms, facilitates one-shot learning of new objects, and improves practical applications. By examining recognition accuracy and understanding human-scale invariance capabilities, we can enhance AI systems and develop more efficient algorithms. Exploring position invariance limitations contributes to our understanding of object recognition factors

and reveals neural computations underlying invariance. Scale and position invariance enable efficient decision-making and object categorization and have applications in security systems and autonomous vehicles. These studies provide insights into cognitive strategies, eye movements, and their role in achieving invariance, benefiting human-machine interfaces and augmented reality systems. Psychophysical investigations into human visual recognition have produced inconsistent and inconclusive results regarding position invariance. While some studies suggest limited position invariance, others propose complete translation invariance for certain shapes. Research on scale invariance with unfamiliar stimuli is scarce, unlike physiological data from monkey studies consistently demonstrating intrinsic invariance within their visual system. Monkeys exhibit scale and translation invariance in the inferior temporal (IT) cortex after becoming familiar with a novel object from a single viewpoint. However, the extent of intrinsic invariant recognition in humans remains uncertain, emphasising the need for further exploration in this intriguing field. The influence of eccentricity on visual acuity in primates has been extensively studied. It has been observed that visual acuity decreases linearly as eccentricity increases. This relationship aligns with the expansion of receptive fields within the primate visual cortex with increasing eccentricity. The concept of the visibility window emerges, indicating the range of visual angles at which objects can be recognised at different sizes. Fine details are perceived in the smaller receptive fields of the foveola, while larger receptive fields at greater eccentricities capture coarser details. These findings deepen our understanding of the intricate relationship between eccentricity, visual acuity, and receptive field sizes in primate visual processing.

This study aims to investigate the existence and characteristics of a distinct region of invariance within the broader region of visibility. It seeks to determine the dimensions and features of this region and explore the correlation between scale, position, and visibility in terms of invariance. The research adopts a one-shot learning approach using unfamiliar stimuli to examine human invariant recognition. Hierarchical convolutional neural networks, specifically eccentricity-dependent neural networks (ENN), are explored as potential models for explaining the experimental data. By combining experiments and simulations, the study aims to unravel neural representations and shed light on the nature of invariant recognition in the brain.

**Methodology:** *(Proper Format)*

Exploration of Perception through Psychophysical Experimentation:

In developing the improvement set, an assortment of 27 Japanese letters were utilized as target protests, every one of which was matched with another Japanese letter filling in as a distractor. This matched course of action should be visible in Figure 1A. Every preliminary comprised of a succession where one of the 27 objective letters was at first introduced as the objective upgrade. Hence, a test letter was shown that could be either indistinguishable from the objective letter or its relating distractor. Calibri letters changed ready and size and were introduced on a uniform white foundation on a Dell U2412M screen with an invigorate pace of 60 Hz. The trial methodology was worked with by the Psychophysics Tool stash 43 for MATLAB 44 running on a Linux PC. To guarantee reliable review conditions, members sat at a proper distance of 1.26 m from the screen and utilized a jaw rest. Scale Invariance Examination: To analyze scale invariance, target and test letters were set in the focal point of the screen while letter size was methodicallly fluctuated. Two blocks of trials were led to evaluate scale invariance in acknowledgment. In the main scale trial block, we tried letter sizes of 30 minutes of circular segment (') and 2 degrees (°). In particular, the objective and

test letter size blends were (30', 30'), (30', 2°), and (2°, 30'). Here, the main component addresses the size of the objective letter, and the subsequent component addresses the size of the test letter. Additionally, in the second block of the scaling test, we utilized letter sizes of 30' and 5°, with the objective and test size blends being (30', 30'), (30', 5°), and (5°, 30°). Similar gathering of subjects partook in both scaled tests, with at least one day between blocks to keep away from any elicitation of the subjects' set improvements. Interpretation Invariance Investigation: To evaluate interpretation invariant acknowledgment, we held steady sizes for both objective and test letters while controlling the place of the test letters comparative with the place of the objective letters. The tried circumstances were partitioned into two gatherings: 1. In the "learning in focal vision" class, members were presented to target letters situated at their visual obsession point, which matched with the focal point of the screen. Inside this condition, the test letters were shown either at a similar situation as the objective letters (alluded to as (0 0)) or at the edge of the members' visual field. This subsequent situation is signified as (0 D), where 0 shows the objective position situated at the focal point of the screen and D addresses the whimsy in visual levels of the place of the test letter comparative with the obsession point. 2. In the "Fringe Vision Learning" class, target letters were introduced to the subject's visual outskirts. Hence, the test letter was introduced in various setups. It could show up at a similar unconventionality as the objective letter, named as (D) or focused (D 0). Furthermore, the test letter can be shown on the contrary side with a similar unconventionality as the objective letter, indicated as (D Opp). In this review, both focal and fringe vision conditions were researched utilizing the accompanying setups: I) whimsy D = 1, 2, 3° with a steady letter size of 30 minutes of circular segment ('); ii) unpredictability D = 2, 2.5° with a letter of size 1 degree (°); and iii) flightiness D = 2, 4, 5, 7° with a letter size of 2°. Bigger letters were utilized to cover a more extensive uprooting range, as the scope of perceivability extends directly with letter size. Because of the bigger number of conditions engaged with interpretation invariance tests contrasted with scale invariance tries, a rehashed set of 27 Japanese letters was utilized in two separate meetings. In the underlying meeting, subjects were given 27 preliminaries and were told to return for the second meeting after a base break of 40 minutes to guarantee that they didn't hold the letters. To manage translational invariance, our exploratory plan included similar gathering of subjects partaking at a few different relocation whimsies for each case. These removal conditions were controlled on various days with an adequate span between them to forestall acquaintance with the upgrades. This restricted redundancy, played out a limit of multiple times, was expected to keep subjects from turning out to be excessively acclimated to the boosts while permitting us to segregate the effect of uprooting from a proportion of invariance unmistakable from individual contrasts between subjects. Two separate gatherings of subjects were engaged with the analyses: one gathering took an interest in the letter-size conditions enduring 30 minutes, while another gathering participated in the 1-grade letter-size conditions. For 2-degree letters, similar gathering of subjects went through testing at D = 2° and 7°, while an alternate gathering was inspected at D = 4° and 5°. It is quite significant that the members engaged with the interpretation invariance tests were unmistakable from those in the scale invariance tests. The improvement set comprised of 27 Japanese letters matched with relating distractor letters, as portrayed in Figure 1A. Every preliminary included the introduction of an objective letter followed by a test letter, which could be equivalent to the objective or its distractor. The Calibri letters were introduced in differing positions and sizes on a Dell U2412M screen with a revive pace of 60 Hz. The trial system used the Psychophysics Tool kit 43 for MATLAB 44 on a Linux PC. Members kept a proper separation of 1.26 m from the screen and utilized a jawline rest for reliable review conditions. Member: To guarantee the curiosity of the improvements and kill any earlier commonality, we painstakingly chose members for the tests who were not familiar with Japanese letters. All

members had typical vision or vision adjusted to ordinary norms. A sum of 10 subjects were selected for the scale-invariance tests, while the interpretation invariance tests required somewhere in the range of 11 and 12 subjects (12 subjects for the 30-minute letter conditions, 11 subjects for the 1-degree letter conditions, and 12 and 11 subjects for the 2-degree letter conditions at D = 2°, 7°, and D = 4°, 5°, separately). If a subject displayed a precision execution underneath 0.6 for the unimportant condition, where the objective and test letters were of a similar size and introduced at the middle ((0 0)), that subject was barred from additional examinations. Essentially, in the event that a subject neglected to meet the standard measures for one uprooting condition, they were likewise prohibited from other relocation conditions. Following the prohibition of subjects underneath the gauge standards, a sum of 10 subjects were remembered for the scale-invariance tests. With respect to the interpretation invariance tests, 9 subjects for every condition were incorporated for the 30-minute letter conditions, 11 subjects for each condition for the 1-degree letter size, and 10 subjects for every condition for the 2-degree letter size. This fastidious choice interaction guaranteed the support of proper subjects for each examination, taking into consideration significant correlations and solid investigations. Also, to confirm the straightforwardness of the planned undertaking and decide the perceivability range for people with earlier information and memory of Japanese letters, we directed tests on three Japanese subjects. It is essential to take note of that for the Japanese members, we utilized a similar trial arrangement and errand. In any case, the goal was not to survey invariant item acknowledgment in a single shot advancing yet rather to assess the perceivability of the letters at various sizes and positions. This examination permitted us to assemble important experiences into the exhibition of subjects with earlier commonality and mastery in Japanese letters while guaranteeing the emphasis stayed on the center exploration question. General Exploratory Strategy: The estimation of letter acknowledgment precision was performed utilizing the equivalent yet unique undertaking. Members were told to focus on a dark spot situated at the focal point of the screen. After a short period, the spot vanished, and an objective letter was momentarily introduced, trailed by time frames white screen. The cycle was rehashed with a test letter. Members then, at that point, demonstrated whether the objective and test letters were something very similar or unique. An example succession of letter introductions should be visible in Figure 1C. Every preliminary comprised of new letter matches, with the test letter arbitrarily picked as indistinguishable or unique in relation to the objective. The brief show time forestalled eye developments and guaranteed seeing at the planned unconventionality. In both scale and interpretation invariance tests, upgrade request were randomized, and the dispersion of the equivalent and various preliminaries and left-right visual field introductions was adjusted. Equivalent quantities of preliminaries were directed for each condition.

**Model Experiments:** The eccentricity-dependent neural network (ENN), shown in Figure 2, leverages two important characteristics of retinal sampling: varying receptive field sizes for specific positions and increasing receptive field sizes with eccentricity. ENN employs weight-sharing and pooling across different positions and scale channels to achieve invariance. In order to assess the model's ability to capture invariant representations of transformations, we conducted comparisons with behavioural data on invariant object recognition.

*Models-* The Eccentricity-dependent Neural Network (ENN), as illustrated in Figure 5, capitalises on two fundamental characteristics of retinal sampling. Firstly, it incorporates receptive fields of varying sizes for specific positions, as demonstrated by previous studies (reference 32). Secondly, it acknowledges that the size of receptive fields increases with eccentricity, as observed in studies (reference 23). Through weight-sharing and pooling

mechanisms across different positions and scale channels, the ENN achieves invariance in its representations. Our hypothesis centred around the model's ability to capture invariant representations of transformations, and thus we conducted tests comparing its performance to behavioural data on invariant object recognition.

The implementation of effective neural networks (ENN) is built upon the foundation of convolutional neural networks (CNN). However, there is a fundamental distinction between ENN and CNN in terms of input processing. In ENN, the input comprises multi-scaled, centred crops of the input images, as depicted in Figure 2B. This approach enables sampling of the central region, corresponding to the foveal area, at multiple resolutions, while the peripheral region is sampled at a lower resolution. Notably, different scale channels share It is worth mentioning that the simulations in this research paper utilised a partially pre-existing implementation.

The ENN model we examined consisted of four layers, including a fully connected layer at the end, resembling the organisation of V1-V2-V4-IT-PFC in the human ventral stream. To establish a correspondence between stimulus size and visual angle, a hyperparameter was introduced where 450 pixels represented 1° of visual angle. This allowed for a more direct comparison between the modelling results and human data. For instance, to extract features of letters with a size of 30′, simulated letters of 225 pixels were placed in the visual field of the model. As mentioned earlier, the input to the model comprised multi-scaled, centred crops of images, with 10 crops used, increasing in size exponentially by a factor of 1.5. The total visual field processed by the model was approximately 19°.
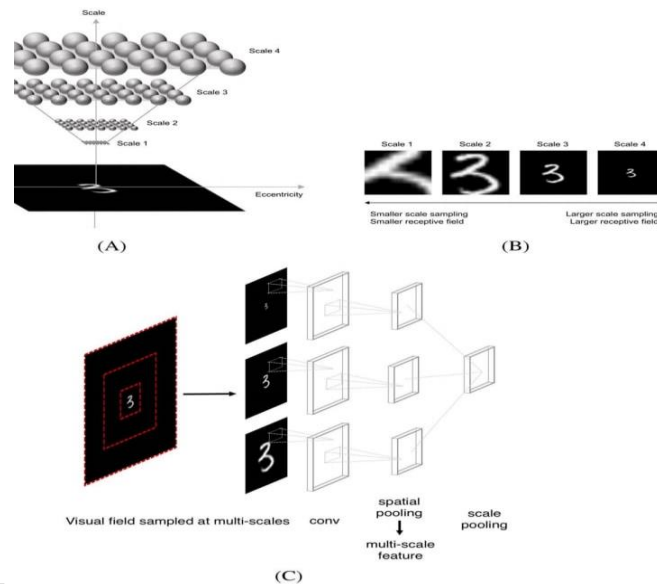
We conducted experiments to evaluate various convolutional and pooling schemes across space and scale, and the results presented here are based on the scheme that exhibited the closest resemblance to human behavioural data. In this scheme, the first layer utilised an 1111-pixel convolutional kernel with a stride of 4 pixels, followed by 555-pixel spatial pooling with a stride of 2 pixels. Subsequent layers employed a convolutional kernel size of 55 pixels with a stride of 1 pixel, along with a pooling kernel size of 55 pixels and a stride of 2 pixels. For scale-pooling, specifically for capturing scale-invariance or extracting features of test letters, 10 scale channels were max-pooled at the final layer.

To ensure the network parameters aligned with human experimental data, we compared the window of visibility for digit recognition. For example, at a distance of 10° from the fovea, humans achieved approximately 67% accuracy for 30' digits. Using linear interpolation, we estimated that accuracy would reach around 77% at 7° for the same digit size. By applying the conversion ratio between pixels and visual angle, we observed a close match between human and ENN accuracy. The conversion ratios and network parameters also aligned with the theoretically estimated size of the smallest receptive fields.

In the case of the Convolutional Neural Network (CNN), the parameters used were identical to those of the ENN, with the exception that it did not involve multi-crop input channels or pooling over scales. Since the CNN model consisted of only one scale channel, the resolution of the input was chosen to match the mid-resolution of the 5th scale channel in ENN.

*Statistical analysis*: The determination of sample sizes in this study was based on the precedent set by previous studies using similar experimental procedures. Statistical methods were not employed for this purpose. The analysis focused on the percentage of correct responses, encompassing both the same and different trials. Parametric tests were used,

assuming a normal data distribution without formal testing. Differences in mean accuracy across multiple conditions were assessed using analysis of variance (ANOVA) or repeated measures ANOVA, depending on the grouping of subjects. Pearson's correlation coefficient (r) was employed to examine feature correlations in simulations.



**Result**- In our study, we examined the intrinsic invariance properties of visual recognition by involving non-Japanese participants in a one-shot learning task with unfamiliar Japanese letters. We presented a target letter followed by a test letter, manipulating their size and placement to investigate the impact of scale and position on recognition. To minimise biases and foveal fixation, letters were randomly displayed in the peripheral visual field, and each letter was presented for a limited duration of 33 ms to ensure consistent visual exposure and prevent eye movements. Diagram 1 provides a visual depiction of the experimental configuration and illustrates the specific set of Japanese letters employed in the study.

## Discussion

This research paper focuses on one-shot learning and investigates the human ability to recognise unfamiliar stimuli despite variations in size and position. The study reveals significant scale invariance in recognition, highlighting the visual system's capacity to generate invariant representations for novel objects. Limited translation invariance is also observed, particularly at lower spatial frequencies. The research provides insights into the intricacies of the human visual system's processing of novel objects and sheds light on factors influencing scale and translation invariance. Comparative analysis reveals the superior performance of effective neural networks (ENNs) over standard convolutional neural networks (CNNs) in capturing and explaining the invariance observed in experimental data. ENNs outperform CNNs in modelling the visual cortex, presenting an alternative computational strategy that leverages multiple small "effective images" at varying resolutions. By recognising objects across resolutions, ENNs achieve scale invariance and offer new insights into visual processing. These findings highlight the potential of ENNs to advance our understanding of the human visual system's mechanisms.

Addressing scale-invariant recognition in convolutional neural networks (CNNs) reveals limitations in their adaptability across datasets. Human perception aligns with effective neural networks (ENNs) enforcing scale-invariant representations. Validating inherent scale invariance highlights the need for meticulous implementation and comparative analysis with neural recordings to refine brain-like computational models.

Future research can explore the impact of effective neural networks (ENNs) on object recognition and eye movements. Investigating the diagnostic role of spatial frequency in ENNs reveals potential scale-invariant critical frequencies and spatial frequency normalisation. However, further investigation is needed to understand the influence of noisy backgrounds and complex imagery on scale channels in ENNs. Additionally, ENNs' positional invariance with low-resolution images suggests a strategy for guiding eye movements, prioritising processing in pivotal positions rather than exhaustive high-resolution analysis. ENNs demonstrate robustness in cluttered environments and the ability to focus on fine details, unaffected by foveal crowding. They excel at detecting small targets amidst complex scenes, leveraging their multi-scale feature extraction capabilities. While not directly comparable to human one-shot learning, augmenting these models with explicit multi-scale sampling enhances their capacity for acquiring knowledge about novel object categories with minimal exemplars. The use of multi-scale channels allows for context-based selection, as suggested by studies conducted by Eckstein et al.

**Conclusion**- The research findings highlight the disparities between convolutional neural networks (CNNs) and effective neural networks (ENNs) regarding their ability to capture scale-invariant perceptions. CNNs tend to develop invariants based on specific instances, while ENNs excel at constructing representations that remain invariant across different scales. This study emphasises the importance of architectural considerations and proposes ENNs as promising models for mimicking the human visual system. Further validation and refinement of the ENN architecture are necessary, along with an investigation into the critical spatial frequencies that contribute to object detection within ENNs. Additionally, the lack of positional changes in ENNs when presented with low-resolution images has implications for understanding eye movements and optimising image processing efficiency. Overall, this research deepens our understanding of the computational strategies involved in achieving scale and translation invariance in human vision.

## References

1. Serre, T., Wolf, L., and Poggio, T. (2005). Object recognition with features induced by the visual cortex IEEE Transactions on Pattern Analysis and Machine Intelligence, 29(3), 411–426,

2. DiCarlo, J. J., & Cox, D. D. (2007). Identification of an ambiguous object Trends in Cognitive Science, 11(8), 333–341.

3. Ulman, S., Vidal-Naquet, M., & Sali, E. (2002). Visual features of intermediate complexity and their use in classification Nature Neuroscience, 5(7), 682–687.

4.  Wallis, G., & Rawls, E.T. (1997): Anomalous face and object recognition in the visual system Advances in Neurobiology, 51(2), 167–194.

5.  Riesenhuber, M., & Poggio, T. (1999). Hierarchical models of object recognition in the cortex Nature Neuroscience, 2(11), 1019–1025.

6.  M. Ahmad, J. F. Al-Amri, A. F. Subahi, S. Khatri, A. Hussain Seh et al., "Healthcare device security assessment through computational methodology," Computer Systems Science and Engineering, vol. 41, no.2, pp. 811–828, 2022.

7.  H. Alyami, M. Nadeem, W. Alosaimi, A. Alharbi, R. Kumar et al., "Analyzing the data of software security life-span: quantum computing era," Intelligent Automation & Soft Computing, vol. 31, no.2, pp. 707–716, 2022.

8.  M. Nadeem, J. F. Al-Amri, A. F. Subahi, A. Hussain Seh, S. Ahmad Khan et al., "Multi-level hesitant fuzzy based model for usable-security assessment," Intelligent Automation & Soft Computing, vol. 31, no.1, pp. 61–82, 2022.

9.  A. Alharbi, W. Alosaimi, H. Alyami, M. Nadeem, M. Faizan et al., "Managing software security risks through an integrated computational method," Intelligent Automation & Soft Computing, vol. 28, no.1, pp. 179–194, 2021.

10. Attaallah, S. Khatri, M. Nadeem, S. Anas Ansar, A. Kumar Pandey et al., "Prediction of covid-19 pandemic spread in kingdom of saudi arabia," Computer Systems Science and Engineering, vol. 37, no.3, pp. 313–329, 2021.

11. F. A. Alzahrani, M. Ahmad, M. Nadeem, R. Kumar and R. Ahmad Khan, "Integrity assessment of medical devices for improving hospital services," Computers, Materials & Continua, vol. 67, no.3, pp. 3619–3633, 2021.

12. Alharbi, M. Faizan, W. Alosaimi, H. Alyami, M. Nadeem et al., "A link analysis algorithm for identification of key hidden services," Computers, Materials & Continua, vol. 68, no.1, pp. 877–886, 2021.

13. Seh, A. H., Ahmad, M., Nadeem, M., Pandey, A. K., Agrawal, A., Kumar, R., & Khan, R. A. (2021). Usable-Security Assessment of Healthcare Software System Through Fuzzy ANP-TOPSIS Method. International Journal of System Dynamics Applications (IJSDA), 10(4), 1-24. http://doi.org/10.4018/IJSDA.304444

14. Alenezi, M., Nadeem, M., Agrawal, A., Kumar, R., & Khan, R. A. (2020). Fuzzy multi criteria decision analysis method for assessing security design tactics for web applications. Int. J. Intell. Eng. Syst, 13(5), 181-196.