

Few-Shot Deep Hashing for Fine-Grained Visual Retrieval Using Self-Supervised Semantic Excavation Networks

Fernando D. Cox

School of Information Technology, University of Cincinnati, Cincinnati, OH, USA.
fernandodcox@uc.edu

Clifford Ray

School of Computing, Clemson University, Clemson, SC, USA.
ray816@clemson.edu

Abstract

The explosive growth of visual data across domains such as biodiversity monitoring, e-commerce, and autonomous systems demands retrieval engines capable of distinguishing subtle inter-class differences among fine-grained categories while operating under severe label scarcity. Few-shot deep hashing has emerged as a promising paradigm to embed images into compact binary codes that preserve semantic similarity and support efficient large-scale search with limited labeled examples. This paper presents a systems-oriented examination of few-shot deep hashing for fine-grained visual retrieval, centering on self-supervised semantic excavation networks that extract rich discriminative structures without exhaustive supervision. We delineate the architectural choices that synthesize self-supervised pretext tasks, asymmetric hash coding, and meta-learning protocols into an integrated framework, and we analyze the structural trade-offs involving code length, retrieval accuracy, computational overhead, and resilience to data shift. Beyond algorithmic design, we investigate the broader system landscape, including cloud-edge deployment topologies, approximate nearest neighbor indexing, resource sustainability, and the socio-technical governance of fairness, privacy, and transparency. Through conceptual modeling and cross-domain illustrations, we argue that robust and equitable visual retrieval cannot be achieved by optimizing hashing objectives alone; it requires a holistic infrastructure that accounts for data curation biases, energy budgets, and regulatory compliance. The article concludes by outlining a forward-looking agenda for responsible few-shot hashing systems, emphasizing the need for interdisciplinary collaboration across machine learning, systems engineering, and policy-making.

Keywords

few-shot learning; deep hashing; fine-grained visual retrieval; self-supervised learning; semantic excavation; system architecture; fairness; sustainability.

1. Introduction

The proliferation of high-resolution image collections has transformed the landscape of visual search, from biodiversity archives cataloging thousands of avian species to product recognition systems that must differentiate nearly identical consumer goods. Fine-grained visual retrieval tasks, where the goal is to distinguish subordinate categories with subtle morphological or textural variations, place extreme demands on both representational fidelity and search efficiency. Traditional deep learning solutions that rely on exhaustive human-

annotated exemplars for every fine-grained class are increasingly untenable due to the labor costs and temporal constraints involved in maintaining up-to-date taxonomies. In response, few-shot learning paradigms have been proposed to imbue retrieval systems with the ability to adapt to novel categories after observing only a handful of labeled instances. Concurrently, the practical deployment of large-scale retrieval services requires drastic compression of visual features into compact binary codes through deep hashing, enabling sub-linear search times and dramatically reduced memory footprints. The confluence of these imperatives motivates the design of few-shot deep hashing architectures that can simultaneously handle label poverty and computational parsimony while preserving the delicate inter-class separability characteristic of fine-grained domains.

This paper adopts a systems-level perspective on the construction of such architectures, with a specific emphasis on self-supervised semantic excavation networks. The core idea behind semantic excavation is to leverage unlabeled data, often abundant in target domains, to autonomously discover latent discriminative structures that can later be exploited when only a few labeled examples are available. By integrating self-supervised pretext tasks with asymmetric hashing constraints, these networks excavate semantic boundaries that are especially valuable for fine-grained differentiation, such as correlations among local part configurations or subtle chromatic gradients. Although numerous algorithmic contributions have been made in deep hashing and few-shot learning individually, the engineering community lacks a comprehensive analysis of how these components interact within a full-stack retrieval system, encompassing training pipelines, indexing infrastructure, edge-cloud orchestration, and the societal dimensions of fairness and accountability. The present work fills this gap by examining the architectural trade-offs, deployment topologies, robustness considerations, sustainability requirements, and governance frameworks that collectively determine the viability of few-shot deep hashing for fine-grained visual retrieval. We do not propose a single novel algorithm; rather, we offer a synthetic conceptual treatment that is intended to guide system architects, policy makers, and interdisciplinary research teams toward more responsible and effective designs.

The exposition is organized as follows. Section 2 reviews the structural evolution of deep hashing systems and situates self-supervised semantic excavation within this trajectory. Section 3 articulates the integration of few-shot learning protocols into hashing architectures and analyzes generalization dynamics under data scarcity. Section 4 explores the deployment infrastructure, indexing strategies, and sustainability concerns that surround large-scale retrieval services. Section 5 addresses robustness, fairness, and governance, proposing a framework for equitable and accountable visual search. Section 6 concludes with a synthesis of insights and a call for interdisciplinary convergence.

2. Foundations and Structural Evolution of Deep Hashing Systems

Early approaches to content-based image retrieval relied on hand-crafted descriptors and quantization techniques that were agnostic to semantic categories. The emergence of learning to hash marked a paradigm shift by optimizing binary codes directly for similarity-preserving properties derived from labeled data [1]. The first generation of deep hashing models coupled end-to-end feature extraction with hash code learning, enforcing pairwise or triplet constraints that pulled semantically similar images together in Hamming space while pushing dissimilar ones apart [5]. These supervised hashing frameworks demonstrated substantial improvements over shallow alternatives but inherited a critical vulnerability: their reliance on abundant

labeled pairs made them ill-suited to settings where annotations are costly or constantly evolving, such as fine-grained species inventories or rapidly changing fashion catalogs.

As the field matured, asymmetric deep hashing formulations gained traction by decoupling the encoding of database items and queries, allowing one stream to operate under relaxed optimization conditions while the other encoded stricter binary constraints [6]. This asymmetry proved effective in scaling to extremely large galleries, because the database side could be optimized with high precision during offline stages while queries benefited from approximate yet fast coding. Nevertheless, supervision remained the primary supervisory signal, and the latent semantic richness of unlabeled data was often discarded. The subsequent shift toward self-supervised learning in computer vision [3] seeded a new class of hashing methods that could exploit large uncurated image sets by designing pretext tasks such as instance discrimination, rotation prediction, or contrastive clustering. In the context of hashing, self-supervision provided a pathway to excavate semantically organized structures without manual labeling, aligning with the broader trend of reducing annotation debt in industrial retrieval pipelines.

The notion of semantic excavation, as operationalized in recent deep hashing research, goes beyond generic self-supervision by explicitly targeting the fine-grained semantic boundaries that are most crucial for subordinate category retrieval. Rather than merely learning transformation-invariant global features, semantic excavation networks construct asymmetric correspondences between augmented views and enforce margin-scalable constraints that adapt the penalties for inter-class confusion based on the estimated semantic proximity of samples [7]. This design encourages the binary codes to capture not only broad category separations but also the nuanced gradients that differentiate closely related fine-grained classes. Critically, such self-supervised excavation can serve as a pre-training phase that creates a semantically organized hash space, which can subsequently be adapted to specific fine-grained tasks with minimal annotation effort. The structural trade-offs in these architectures revolve around the tension between the richness of the excavated semantic fabric and the compactness of the hash codes. More aggressive excavation may produce highly discriminative yet fragile representations that overfit to the self-supervised task distribution, whereas underexcavated codes may fail to separate fine-grained categories when only a few examples are provided.

HashNet and related continuation methods introduced the insight that learning discrete binary codes can be stabilized by gradually annealing the approximation of the sign function during training [8]. When combined with self-supervised objectives, such continuation strategies enable the network to smoothly traverse the loss landscape from a relaxed continuous embedding space to a highly compressed binary manifold without catastrophic semantic collapse. Furthermore, adversarial self-supervised hashing models have demonstrated that incorporating a generative adversary can improve the uniformity of hash code distributions, an important property for maximizing the information capacity of compact codes [9]. These architectural innovations collectively illustrate a design space in which the choice of pretext task, the symmetry of encoding pathways, the hardness of binary approximation, and the method of semantic boundary refinement must be co-optimized to suit the domain’s fine-grained granularity and the anticipated few-shot scenarios. A system architect must therefore evaluate not only the retrieval accuracy on a fixed benchmark but also the resilience of the excavated semantics under distribution shifts, the training stability at scale, and the compatibility with downstream few-shot adaptation modules.

3. Few-Shot Adaptation and Self-Supervised Semantic Excavation Architectures

Integrating few-shot learning into deep hashing architectures requires rethinking the training protocol and the internal representational geometry. In standard metric-based few-shot learning, a model is meta-trained across many episodic tasks, each involving a small support set and a query set drawn from a base set of classes, so that it learns to rapidly generalize to novel classes [4]. Prototypical networks, which compute a class prototype as the mean embedding of the support examples, have been especially influential. When the embedding function maps to a binary Hamming space rather than a continuous Euclidean space, the prototype computation becomes nontrivial, because the binary constraint makes simple averaging ambiguous. Recent meta-hashing frameworks address this by learning to fold few-shot adaptation into the hash coding process itself, for example by employing attention mechanisms that weigh the support examples' binary codes according to their relevance to a query [10]. The self-supervised semantic excavation stage becomes a powerful precursor for such meta-hashing, because the pretext tasks endow the hash space with a rich topology that already respects fine-grained semantic clusters even before any labeled support set is provided.

The interplay between self-supervised excavation and few-shot adaptation can be understood as a two-phase knowledge transfer pipeline. In the first phase, the network ingests vast amounts of unlabeled data and excavates a latent semantic manifold by solving asymmetric instance-level discrimination tasks, while simultaneously enforcing margin-scalable constraints that sharpen the boundaries between visually similar yet semantically distinct regions [7]. The resulting hash space exhibits a property that might be termed semantic prototypicality: visually coherent groups tend to form tight clusters with well-separated margins, even though no explicit labels were used to define these clusters. In the second phase, when a novel fine-grained category is introduced with only a handful of labeled examples, the meta-learned adaptation module can rapidly identify the appropriate region of the hash space and refine the decision boundaries by re-weighting the geometric relationships that the excavation phase established. This two-phase regime is structurally analogous to the pretrain-then-fine-tune paradigm of natural language processing, but it operates under the additional constraint of extreme compression.

Several system-level trade-offs emerge from this design. The capacity of the self-supervised excavation module, typically implemented as a large backbone network such as a deep residual architecture [21], directly influences both the quality of the excavated semantics and the computational footprint of the system. Deploying heavy excavation on resource-constrained edge devices is impractical, so a common strategy is to perform excavation offline in the cloud and distribute a frozen hash encoder to edge nodes. However, this division introduces a distributional fragility: if the unlabeled data used for excavation does not adequately represent the target deployment environment, the few-shot adaptation phase may encounter a domain gap that degrades retrieval performance. Thus, architectural decisions about the granularity of domain-specific unlabeled pre-training and the frequency of encoder updates become tightly coupled with deployment topology.

Another critical consideration is the code length used for hashing. Fine-grained retrieval demands longer binary codes to preserve enough discriminative information to separate highly similar categories, yet longer codes increase storage overhead and slow down Hamming distance computations at query time. Self-supervised semantic excavation can partly mitigate this tension by concentrating the code's representational budget on the most semantically salient dimensions, effectively increasing the per-bit information content.

Margin-scalable constraints play a key role here: by imposing larger margins for categories that are visually confusable, the excavation process allocates more coding capacity to the most challenging distinctions. When combined with few-shot adaptation, this allocation becomes even more critical, because the limited support samples cannot independently correct for an overly compressed hash space that collapsed fine-grained boundaries. System designers must therefore calibrate excavation objectives and code lengths through rigorous cross-validation that simulates realistic few-shot scenarios, rather than relying on standard classification metrics alone.

Robustness to label noise and open-set conditions further complicates the few-shot hashing pipeline. In many real-world fine-grained domains, the very definition of a category may be fluid, and support sets may contain mislabeled images. Self-supervised excavation provides a degree of label-noise immunity by anchoring the hash space in unlabeled distributional properties rather than brittle human annotations. Nevertheless, the adaptation phase remains vulnerable to corrupted support sets, a vulnerability that can be partially alleviated by incorporating robust prototype estimation techniques, such as iterative re-weighting of support examples based on their consistency with the pre-existing semantic clusters. The governance implications of these robustness properties will be revisited in Section 5.

4. Infrastructure and Deployment for Large-Scale Fine-Grained Retrieval

The operational deployment of a few-shot deep hashing system for fine-grained retrieval demands a carefully orchestrated infrastructure that spans data ingestion, offline model training, hash code generation, index construction, and online query serving. The self-supervised excavation phase is typically the most computationally intensive component, requiring hundreds of GPU-hours to process large unlabeled corpora and converge the excavation objectives. This phase is best executed in a centralized cloud environment with elastic compute resources, where distributed training strategies can be employed to handle the massive data volumes. The resulting semantic hash encoder can then be compressed, quantized, and pushed to edge inference nodes or embedded directly into mobile applications for on-device hash code extraction. The decoupling of training and inference is a standard pattern, but it introduces governance challenges around model versioning, update propagation, and consistency of the hash space across distributed nodes.

Once hash codes are generated for the gallery database, they must be indexed for efficient Hamming distance search. Modern approximate nearest neighbor libraries, such as FAISS, provide highly optimized primitives for billion-scale binary search on GPUs and CPUs [11]. Multi-index hashing strategies further accelerate retrieval by splitting the binary code into multiple substrings and building separate tables for each substring, trading a small amount of recall for significant speed gains. Product quantization offers an alternative approach that can be applied to the continuous feature representations before binarization, thereby reducing the distortion introduced by extreme compression [20]. The choice of indexing structure must account for the unique characteristics of fine-grained retrieval: the distance margins between correct and incorrect matches are often minuscule, so the indexing scheme must preserve Hamming distance ordering with minimal approximation error. Self-supervised excavation, by increasing the inter-class separability in the binary space, can make the retrieval system more tolerant of the approximation losses inherent in multi-index or cell-probe methods, thereby allowing faster query latencies without catastrophic deterioration of precision.

Sustainability considerations permeate every layer of this infrastructure. Training large self-supervised models leaves a non-negligible carbon footprint, and the practice of periodic re-

excavation to incorporate new unlabeled data can quickly accumulate energy costs [16]. System designers must weigh the environmental impact against the retrieval efficiency gains that compact hashing provides at query time. A promising direction is the adoption of progressive excavation strategies, where only partial re-training is performed on newly acquired data, combined with online clustering techniques that update hash assignments incrementally. Furthermore, the energy efficiency of the inference stage can be optimized by selecting backbone architectures that maximize accuracy per FLOP, using neural architecture search to discover lightweight encoders that still preserve the fine-grained semantics necessary for reliable few-shot adaptation. The trade-off between model complexity and retrieval quality must be evaluated not only in terms of mean average precision but also in terms of energy-delay product, a metric that captures the operational cost of serving queries under latency constraints.

The edge-cloud continuum introduces additional architectural complexities. In biodiversity monitoring applications, for instance, field devices with limited connectivity may capture images of rare species and need to perform on-device retrieval against a local hash index that is periodically synchronized with a cloud repository. Fine-grained hashing models deployed on such devices must be capable of few-shot adaptation using only the on-device support set, without access to the original unlabeled corpus. This requires that the self-supervised excavation be comprehensive enough to generalize to the edge environment, yet the hash encoder must be small enough to fit within the memory and computational budgets of embedded hardware. Federated learning frameworks, where multiple edge devices collaboratively update a shared hash space without sharing raw images, offer a potential solution, but they introduce new challenges related to the heterogeneity of edge data distributions and the preservation of fine-grained cluster integrity under federated aggregation. These system-level design tensions underscore that few-shot deep hashing cannot be reduced to a purely algorithmic problem; it is fundamentally an infrastructure design problem that spans hardware, networking, and data management.

5. Robustness, Fairness, and Governance Framework

The deployment of few-shot deep hashing for fine-grained visual retrieval in socially consequential domains such as law enforcement, medical diagnosis, or biodiversity conservation necessitates a thorough examination of robustness, fairness, and governance. Robustness refers to the system’s ability to maintain reliable retrieval performance under distributional shifts, adversarial perturbations, and noise in the few-shot support samples. Self-supervised semantic excavation inherently provides a degree of robustness by anchoring the hash space in unlabeled data patterns that are often more stable than human-assigned labels [7]. However, adversaries can still craft imperceptible perturbations that cause a query image to map to a distant hash code, thereby evading retrieval or fetching incorrect fine-grained results. Hashing systems are particularly exposed to such attacks because the binary quantization step amplifies small input perturbations into large code alterations. Defensive strategies, including adversarial training during the excavation phase and randomized quantization at inference time, must be integrated into the system architecture to harden the retrieval service against adversarial exploitation.

Fairness in visual retrieval is a multi-faceted concern that extends beyond conventional classification bias. In fine-grained domains, the very taxonomy used to define categories may reflect historical inequities or incomplete geographic sampling. For instance, a bird species retrieval system trained predominantly on images from North American and European

sources may exhibit substantially degraded performance for subspecies endemic to under-resourced regions, not because the visual features are intrinsically harder but because the self-supervised excavation phase was starved of representative unlabeled data. This distributional unfairness can be amplified in few-shot scenarios, where the support examples for rare categories are themselves drawn from the skewed distribution. The resulting retrieval system may systematically fail to surface relevant results for certain geographic or demographic user groups, perpetuating a cycle of invisibility. Fairness-aware hashing therefore demands deliberate data curation policies that ensure geographic and demographic representativeness in both the unlabeled excavation corpus and the few-shot support sets [15]. Technical interventions such as adversarial debiasing or fairness constraints on the Hamming space can partially mitigate these biases, but they cannot substitute for foundational governance decisions regarding data sourcing and inclusion.

The governance architecture must also address privacy and transparency. Hash codes, despite their compactness, can leak information about the input image or the identity of individuals, especially when coupled with trained decoder networks. Deployments in sensitive settings therefore benefit from differential privacy mechanisms that inject calibrated noise into the hash code computation, ensuring that the presence or absence of a specific image in the gallery cannot be confidently inferred [18]. However, the introduction of differential privacy lowers the effective bit-rate of the hash code and may degrade fine-grained discriminability, creating a direct trade-off between privacy protection and retrieval accuracy. Transparent documentation of these trade-offs, for example through model cards that specify the privacy budget, the demographic composition of training data, and the expected retrieval performance across subgroups, is an essential component of responsible deployment. Regulatory frameworks such as the General Data Protection Regulation impose additional requirements on the right to explanation, compelling system operators to justify why a particular fine-grained result was retrieved. In the context of deep hashing, achieving interpretable retrieval may involve generating saliency maps that highlight which visual regions most influenced the hash code similarity, a capability that is still nascent in binary embedding research.

Policy implications extend to the governance of AI-powered biodiversity monitoring and surveillance. Self-supervised semantic excavation networks can be repurposed for dual-use applications, tracking endangered species for conservation while also enabling illicit surveillance of human populations in protected areas. The international policy community must therefore develop guidelines that distinguish between benevolent fine-grained retrieval and invasive monitoring, and system designers must build access controls and audit trails into the retrieval infrastructure. The principle of data minimization, central to modern privacy regulation, aligns naturally with the ethos of hashing, since binary codes discard substantial visual detail while retaining semantic gist. Yet this alignment should not be taken for granted; systematic red-teaming is required to assess whether the excavated semantics inadvertently preserve protected attributes or enable re-identification.

Establishing a fairness and governance framework for few-shot deep hashing involves technical, organizational, and legal layers. At the technical layer, audit tools that measure performance disparities across user-defined subgroups must become a standard component of the deployment pipeline, much like unit tests in software engineering. At the organizational layer, interdisciplinary stewardship committees, including domain experts in fine-grained taxonomy, ethicists, and representatives of affected communities, should oversee the curation of unlabeled data and the specification of few-shot tasks. At the legal layer, compliance with

evolving AI regulations necessitates continuous monitoring of retrieval fairness metrics and the maintenance of records that demonstrate due diligence. By weaving these layers together, system architects can move beyond a narrow optimization mindset and toward a socio-technical conception of retrieval quality that accounts for accuracy, equity, transparency, and accountability simultaneously [17, 19]. The self-supervised semantic excavation paradigm, with its dependence on uncurated data and its opaque learned semantics, is particularly in need of such enveloping governance safeguards.

6. Conclusion

Few-shot deep hashing grounded in self-supervised semantic excavation networks represents a compelling systems solution for fine-grained visual retrieval, merging the efficiency of binary coding with the flexibility of learning from limited labels. This paper has argued that realizing the full potential of this paradigm requires moving beyond isolated algorithmic innovations and embracing a holistic, infrastructure-oriented perspective. The architectural choices that shape semantic excavation—the symmetry of encoding pathways, the margin-scalable constraints, the pretext tasks, and the integration of meta-learning for few-shot adaptation—must be evaluated not only on benchmark accuracy but also on their robustness to distributional shifts, their computational and energy costs, and their compatibility with privacy-preserving mechanisms. Deployment infrastructures that span cloud training, edge inference, and distributed index structures introduce further trade-offs between latency, storage, and retrieval fidelity, while sustainability concerns urge caution against indiscriminate large-scale re-training.

Equally critical are the fairness and governance dimensions. Fine-grained retrieval systems can perpetuate representational harms if the self-supervised excavation phase is trained on unrepresentative data, and few-shot settings can amplify these inequities when support sets are themselves biased. Robust governance frameworks that mandate transparency, fairness auditing, differential privacy, and community oversight are not optional add-ons but constitutive elements of a responsible retrieval service. As visual search becomes embedded in public infrastructure—from wildlife monitoring networks to consumer platforms—the integration of technical performance with ethical accountability will determine the societal legitimacy of these technologies. We therefore call for sustained interdisciplinary collaboration among machine learning researchers, systems engineers, domain scientists, ethicists, and policymakers to chart a trajectory in which few-shot deep hashing serves the public good while respecting the complex fabric of our visual world.

References

1. Wang, J., Liu, W., Kumar, S., & Chang, S. F. (2016). Learning to hash for indexing big data—A survey. *Proceedings of the IEEE*, 104(1), 34–57.
2. Lin, T. Y., RoyChowdhury, A., & Maji, S. (2015). Bilinear CNN models for fine-grained visual recognition. *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 1449–1457.
3. Chen, T., Kornblith, S., Norouzi, M., & Hinton, G. (2020). A simple framework for contrastive learning of visual representations. *Proceedings of the International Conference on Machine Learning (ICML)*, 1597–1607.
4. Snell, J., Swersky, K., & Zemel, R. S. (2017). Prototypical networks for few-shot learning. *Advances in Neural Information Processing Systems (NeurIPS)*, 30, 4077–4087.

5. Li, W. J., Wang, S., & Kang, W. C. (2016). Feature learning based deep supervised hashing with pairwise labels. *Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI)*, 1711–1717.
6. Jiang, Q. Y., & Li, W. J. (2018). Asymmetric deep supervised hashing. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1), 3342–3349.
7. Yu, Z., Wu, S., Dou, Z., & Bakker, E. M. (2022). Deep hashing with self-supervised asymmetric semantic excavation and margin-scalable constraint. *Neurocomputing*, 483, 87–104.
8. Cao, Z., Long, M., Wang, J., & Yu, P. S. (2017). HashNet: Deep learning to hash by continuation. *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 5608–5617.
9. Zhan, C., Yu, J., & Tao, D. (2019). Self-supervised adversarial hashing networks for cross-modal retrieval. *IEEE Transactions on Image Processing*, 29, 1800–1813.
10. Cheng, X., Li, X., Yang, Y., & Hauptmann, A. G. (2021). Meta-hashing for large-scale image retrieval. *IEEE Transactions on Image Processing*, 30, 4958–4971.
11. Johnson, J., Douze, M., & Jégou, H. (2019). Billion-scale similarity search with GPUs. *IEEE Transactions on Big Data*, 7(3), 535–547.
12. Norouzi, M., Fleet, D. J., & Salakhutdinov, R. R. (2012). Hamming distance metric learning. *Advances in Neural Information Processing Systems (NeurIPS)*, 25, 1061–1069.
13. Qian, X., Tang, Y. Y., Yan, Z., & Huang, K. (2019). Deep binary representation learning for fine-grained image retrieval. *IEEE Transactions on Image Processing*, 28(10), 5052–5064.
14. Singh, A., & Joachims, T. (2018). Fairness of exposure in rankings. *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2219–2228.
15. Buolamwini, J., & Gebru, T. (2018). Gender shades: Intersectional accuracy disparities in commercial gender classification. *Proceedings of the 1st Conference on Fairness, Accountability and Transparency*, 77–91.
16. Lacoste, A., Luccioni, A., Schmidt, V., & Dandres, T. (2019). Quantifying the carbon emissions of machine learning. *arXiv preprint arXiv:1910.09700*.
17. Mittelstadt, B. D., Allo, P., Taddeo, M., Wachter, S., & Floridi, L. (2016). The ethics of algorithms: Mapping the debate. *Big Data & Society*, 3(2), 2053951716679679.
18. Dwork, C., & Roth, A. (2014). The algorithmic foundations of differential privacy. *Foundations and Trends in Theoretical Computer Science*, 9(3–4), 211–407.
19. Floridi, L., Cowsls, J., Beltrametti, M., Chatila, R., Chazerand, P., Dignum, V., ... & Vayena, E. (2018). AI4People—An ethical framework for a good AI society. *Minds and Machines*, 28(4), 689–707.
20. Jégou, H., Douze, M., & Schmid, C. (2011). Product quantization for nearest neighbor search. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(1), 117–128.

21. He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 770–778.
22. Liu, J., Zhang, Q., & Zhao, Y. (2020). Edge intelligence: Paving the last mile of artificial intelligence with edge computing. *Proceedings of the IEEE*, 108(10), 1781–1802.