

Toward a Computational Model of Transfer

Daniel Oblinger

■ *The Defense Advanced Research Projects Agency (DARPA) explored the application of transfer — a notion well studied in psychology — to machine learning. This article discusses the formal measure of transfer and how it evolved. We discuss lessons learned, progress made at the formal and algorithmic levels, and thoughts about current and future prospects for the practical application of this technology.*

The Defense Advanced Research Projects Agency (DARPA) Transfer Learning Program (TLP) explored the application of “transfer”—a notion well studied in psychology—to machine learning, where it was still novel. The aims of TLP were to understand and formally frame how this intuitively compelling psychological idea might apply in the computational context, build computational models of transfer learning (TL), and explore how these models might apply to practical learning tasks. TLP and the field as a whole made great strides in each of these dimensions. Indeed, the program has helped TL become a recognized subdiscipline of machine learning. Other articles in this special issue detail the work accomplished in TLP; this article focuses on a broad framing of the research conducted and an assessment of its progress, limitations, and challenges, from an admittedly personal but DARPA-influenced perspective.

Framing Transfer Learning

Traditionally every DARPA program has focused its research by requiring a precise measure of progress. The DARPA TLP decided to measure transfer by comparing the learning of tasks A and B versus the learning of B alone. In figure 1 the curve labeled B represents a traditional learning curve of the performance on target task B as a function of the number of training instances. Curve $A + B$ represents the same learning algorithm, given the same sequence of training instances for task B but additionally provided all available training data from task A (suitably transformed) prior to receiving any training data for task B . Intuitively the area between the two curves represents the transfer or “boost” that the algorithm received from exposure to source task A 's data.

TLP developed multiple measures of this area to assess transfer but observed more than a dozen forms of degenerate behavior while trying to apply these seemingly straightforward metrics. For example, some learning algorithms were wildly nonmonotonic, making it hard to select an appropriate cutoff window. That choice affected transfer scores as well as the meaning and relative importance of transfer metrics, such as the jump-start and asymptotic advantages (figure 1). Once these difficulties were mitigated, the primary remaining shortcoming of the metric was that it provided no way to calibrate observed transfer against a baseline (analogous to comparing induction to the a priori most likely class) or any kind of upper limit on obtainable transfer. Yet this basic characterization of transfer was critical to the success of TLP as it afforded an intuitive and practical way to measure progress quantitatively with respect to specific problem classes. Without this measure, documenting research progress would have been impractical from year to year—and quantitative measures are often the most compelling way to justify long-term research investment in an area.

Characterizing transfer by its boost also allowed us to frame the most salient distinctions between the types of transfer explored within TLP. Some algorithms expected to receive their boost as a jump start—an increased y-intercept value (figure 1). Of course, one expects this increased value to be carried across the learning curve, resulting in an expanded area between the curves. Other algorithms promised a greater slope for the $A + B$ curve over some indeterminate “productive” range of the learning algorithm. On inspection, the first class of learning algorithms often attempted to use heuristics to map learned knowledge directly from instances of the source task A to instances of the target task B , while the latter attempted to map some form of bias (for example, expressed as joint Bayesian priors on the source and target learning tasks).

Interestingly the different approaches had divergent but harder to characterize differences in asymptotic performance as well, which is often of the greatest import for applications of interest to DARPA. Transfer from instances often either had no asymptotic difference or an extreme asymptotic difference in cases where the learning algorithm failed without the generalized knowledge transferred from those instances. Transfer of bias or priors tended to provide modest but consistent differences in asymptotic performance. TLP never found a way to predict these effects.

Another key, though perhaps expected, difference between these approaches is how each would naturally frame the transfer problem itself. The transfer of bias typically yields the largest improvements when employed between large families of

related learning tasks, whereas direct mapping of knowledge applies more naturally to the A -to- B transfer case addressed within TLP.

Status and Prospects for a Formal Understanding of Transfer

A theoretical understanding of transfer is still in its infancy. By the end of TLP, we had a relatively intuitive, formally characterized answer to the question, “How much transfer has occurred between these two learning tasks?” But TL is far from having a theoretical basis the equivalent of computational learning theory, in which the Vapnik-Chervonenkis dimension provides both a formally pleasing and practically useful measure of inductive difficulty. Creating a formal theory of transfer remains a critical, yet difficult, direction for future work.

Although there were some attempts, we did not arrive at practical definitions for many of the concepts that we nonetheless treated as meaningful throughout TLP. Notions of “distance” between different inductive learning tasks, “difficulty” of transfer, and “types” of transfer seemingly could be formally characterized (with sufficient constraints). We assumed these notions were real, but left them underspecified. Definitions are needed that are intuitive, provide explanatory power, and are practically measurable.

Today we use the term *transfer* to cover a number of qualitatively distinct processes for connecting related learning tasks. Even making coarse-grained qualitative statements that apply to all transfer algorithms is difficult. Despite all these shortcomings, TLP produced algorithms that provide significant performance improvements across many practical transfer tasks.

Status and Prospects for Research on Algorithms and Applications

Over the course of TLP, we applied transfer to dozens of application tasks ranging from simple synthetic tasks designed to explore specific classes of transfer, to complex groups of learning tasks in real-world sensing and acting. In this section we describe the two transfer tasks used in TLP’s final phase and use them to consider the status, challenges, and opportunities of the two types of transfer attempted.

Cognitive Approaches

Cognitive approaches (broadly, those that explicitly map knowledge from instances of the source task to instances of the target task) require a rich representational space for tasks in both domains. Cognitive approaches were applied to the task of clas-

sifying and schematically describing football plays based on overhead camera footage of college football games. Knowledge gained from learning in this domain was then transferred to the task of actually *playing* football in a computer simulation.

There are profound differences between these source and target learning tasks. Beyond the obvious differences between classifying a video stream and playing a computer game, the rules (even the number of players) varied between these tasks. Thus the only meaningful transfer concerned the coordinated behaviors of implicitly defined groups of players, for example, strategies involving causal relationships among the center, quarterback, and receiver. Given the enormous parametric space of possible strategies for coordinating players, without transferred knowledge the learning algorithms were hopelessly lost in searching for even the simplest playing strategies (for example, having the quarterback drop back at the beginning of a play, having receivers run laterally to avoid defensive coverage). Causal explanations for observed play behaviors provided an extremely powerful bias, guiding reinforcement learning to far more complex playing strategies. This highlights both the potential and the constraints on practical application of these cognitive approaches. The benefits will almost always be transformative for the target learning task, but they will only come in a context with deep and rich task knowledge tying the two transfer domains together.

Such task knowledge is not easy to encode and does not exist for most practical applications. One approach attempts to generalize the learning method such that the knowledge mapping and acquisition required for transfer utilize the same common underlying mechanism. Another cognitive approach, pursued by many researchers, treats transfer as one of multiple distinct forms of mental processing. The key challenge to widescale application of this technology is creation of mechanisms for capturing the background knowledge.

Bayesian Approaches

In TLP the Bayesian approaches

involved engineering research to characterize the space of relationships between source and target as a set of joint variables over which the source task can provide priors used to constrain the target task. Framed in this way, each new class of transfer tasks becomes an exercise in encoding expected relations as a Bayesian inference problem. The final TLP task here was to transfer the *recognition* of objects (and in some cases parametric models of objects as a composition of primitive parts) from still images to the robotic *manipulation* of related physical objects. Specifically the parametric model of recognized object classes was used to define the priors on grasping strategies for physical objects.

Transfer of priors is now being used in other application domains. DARPA is funding work on text extraction (for example, extraction of relevant parameters on natural disasters as reported in news articles) that transfers specific forms of linguistic knowledge between different languages, genres (for example, AP articles versus Twitter traffic), and extraction domains (for example, natural disasters, crime reporting, mergers and acquisitions). The mechanism for communicating the relationships between learning tasks as priors on those tasks appears to be quite general and practical. It requires a sophisticated understanding of both statistical domain modeling and the targeted transfer; however, the results have dramatically reduced the training required.

A key challenge for Bayesian approaches is to identify a theory that enables constructing generic building blocks that would allow new classes of transfer to be cost-effectively expressed as a combination of previously characterized transfer components.

The Future of Cognitive and Bayesian Approaches

Transfer learning performed by either a cognitive or Bayesian approach could, at least in principle, be encoded within the other approach. In practice, however, the two research areas appear to be moving in quite different directions. Perhaps the most important difference lies in their assumptions about the

source of the knowledge underlying transfer. In the Bayesian approach it is natural to assume that engineers hand-code the underlying knowledge as priors and joint variables, based on understanding the relationships between learning tasks. In the cognitive approach the form of knowledge is more amenable to automated acquisition. The Transfer Learning Program didn't address the acquisition of underlying knowledge, but the two approaches have very different objectives. The Bayesian approach anticipates a toolkit in an emerging engineering discipline of transfer; a mature, successful cognitive approach would provide transfer as part of an autonomous learning agent.

The two approaches also contrast sharply in the complexity of their underlying connective knowledge. Bayesian approaches generally start with an extremely small bit of knowledge relating source and target tasks (often a single organizing principle), using this one insight to structure the joint variables and priors for transfer. By contrast, cognitive approaches encode a much deeper theory connecting the learning tasks. Not surprisingly, in TLP the cognitive approaches often needed far less training data to achieve good transfer performance, while the Bayesian approaches required less engineering time to encode their one principle. Because of the limited development required, specific instances of the Bayesian approach to transfer are probably closer to practical application today; the cognitive approaches are aiming toward general models of intelligence.

Daniel Oblinger is focused on scalable approaches for knowledge acquisition, typically applying inductive learning, but in rich environments where time, interaction, and structure make the traditional "vector of features" model ineffective. He currently serves as a program manager at DARPA running the Machine Reading and Bootstrapped Learning programs. Prior to this he was a research staff member at IBM T.J. Watson Labs where he developed inductive algorithms appropriate for a range of tasks from predicting speech and reading deficiencies in programming by demonstration.