

# Computational Modeling of Variability in the Balance Scale Task

Fiona M. Richardson (f.richardson@bbk.ac.uk)

Frank D. Baughman (f.baughman@bbk.ac.uk)

Neil Forrester (n.forrester@bbk.ac.uk)

Michael S.C. Thomas (m.thomas@bbk.ac.uk)

Developmental Neurocognition Laboratory,  
School of Psychology, Birkbeck College,  
University of London, WC1E 7HX UK

## Abstract

The study of variability in reasoning can shed light on several issues, including mechanisms underlying developmental change, individual differences, and developmental disorders. We explored the basis of variability in a much-studied task in cognitive development, the balance scale. Starting with a simple feed-forward connectionist model and training patterns based on McClelland (1989), we investigated computational parameters, problem encodings, and training environments that contribute to variability in development, both across groups and within individuals. We report on the parameters that affect the complexity of reasoning and the nature of ‘rule’ transitions exhibited by networks learning to reason about balance scale problems.

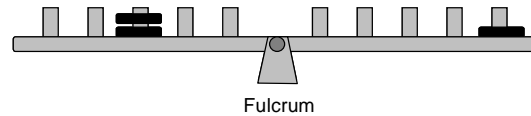
## Introduction

In the study of cognitive development, a rich literature has accumulated on the balance scale task. In this task, different numbers of weights are placed at distances either side of a fulcrum and the child is asked whether the scale will balance, tip left, or tip right when released (Inhelder & Piaget, 1958). For both recent empirical and computational approaches to this domain, the cornerstone is Siegler’s initial work (1976, 1981) in which children’s decisions at different ages were characterized in terms of four rules of increasing complexity. Rules *I* to *IV* describe the child’s performance on each of the six different problem types (see Figure 1), with *Rule IV* representing mastery.

Siegler’s rule assessment methodology has provoked much debate, both with regard to whether the rules he postulated are sufficient to capture children’s behavior (e.g., Wilkening & Andersen, 1982) and whether rules actually play a causal role in driving behavior (Hardiman, Pollatsek, & Well, 1986). Rule-based theories of development have traditionally struggled to explain the mechanisms mediating transitions between rule states, leading to theories based on connectionist learning models (e.g., McClelland, 1989). As an example of the debate, it has been argued that children use different rules depending on the torque value (where torque = weight x distance from fulcrum): Ferretti and Butterfield (1986) found that

problems with a large difference in the torque acting on each side were likely to draw responses consistent with a more advanced rule; Jansen and van der Maas (1997) later reported that the *torque difference* effect only occurred for problems with extreme torque values.

|   |   |
|---|---|
| 1. <u>B</u> alance problems   | The same weight is positioned at the same distance on both sides            |
| 2. <u>W</u> eight problems  | Different weights are placed at the same distance on both sides             |
| 3. <u>D</u> istance problems  | The same weight is placed at different distances either side of the fulcrum |
| Conflict problems: <i>weight and distance are placed in conflict on the two sides</i> |   |
| 4. <u>C</u> onflict <u>W</u> eight  | The side with the larger weight tips  |
| 5. <u>C</u> onflict <u>D</u> istance  | The side with the weights the largest distance from the fulcrum tips        |
| 6. <u>C</u> onflict <u>B</u> alance   | The scale balances  |



|           |  |
|-----------|--|
| Rule I    | Consider weight only; the side with the most weight tips   |
| SDD       | Smallest Distance Down: If distance differs, select the side with the weights closest to the fulcrum   |
| Rule II   | Consider the distance dimension if weights on either side are equal  |
| QP        | Qualitative Proportionality: For conflict problems scale will balance, as larger weight on one side will compensate for greater distance on other  |
| Rule III  | Consider information on weight and distance but (as unable to combine them) guess on conflict problems   |
| Rule IIIa | Focus on either the weight or distance and make a perceptual decision on which side will tip   |
| DD        | Distant Dominant: side with weights largest distance from the fulcrum tips   |
| Addition  | Calculate ( <i>weight + distance</i> ) on each side of the fulcrum; side with highest value tips   |
| Buggy     | For side X with more weights but smaller distance, shift weights away from fulcrum until the distance on each side is equal; for each shift, remove one weight from X. Side with greater final weight tips |
| Rule IV   | Mastery: solve problems by calculating the torque ( <i>weight x distance</i> ) on each side of the fulcrum; side with highest torque value tips  |

Figure 1. The balance scale and developmental ‘rules’

Despite criticism, Siegler’s rules have stood the test of time, albeit with proposed additions (and replacements) to the original four core rules. For example, the *smallest distance down* rule (*SDD*, Figure 1) has been proposed as a rule used by children only when in transition between rules *I* and *II* (Jansen & van der Maas, 2002). The majority of new rules have emerged through the scrutiny of behavior surrounding *Rule III* where, according to Siegler’s scheme, children perform well when either weight or distance

information unambiguously predicts the side to tip, but then guess when these sources of information conflict. Some of the new rules proposed to account for the variability around *Rule III* include: *Rule IIIa*, the *qualitative proportionality*, *distance dominant*, *addition*, and *buggy* rules (Ferretti & Butterfield, 1986; Jansen & van der Maas, 1997, 2002; Normandeau et al., 1989; van Maanen, Bean & Sijtsma, 1989; Wilkening & Andersen, 1982).

The existence of additional rules has found support from Latent Class Analysis, a statistical technique for categorizing behavioral data into consistent subgroups (e.g., Jansen & van der Maas, 1997, 2002). Though these analyses differ in the number of classes generated (relating to a free parameter in the technique), they converge on the idea that *Rule III* behavior consists of a variety of strategies that children tend to switch between. Recent work examining reaction times (RT) as well as accuracy has supported the development of more complex balance-scale strategies with age, favoring the *buggy* rule over the *addition* rule as a *Rule III* strategy (van der Maas & Jansen, 2003), although the response patterns for *buggy* and *addition* are equivalent.

Individual variability in performance on different problem types has been acknowledged in theories of the phases of development. The *staircase model* captures the phases of development by proposing that transitions between rules are quick with relatively little overlap, while transitions in the *overlapping waves model* are more gradual and interleaved, particularly around *Rule III* (Siegler, 2002). A combination of these two models (Jansen & van der Maas, 2002) captures the behavioral data via steep transitions between *Rule I* and *Rule II* but overlap and gradual transitions between subsequent rules (such as *Rule II*, *Rule III*, and the *addition* rule) prior to reaching *Rule IV*.

Computational approaches have sought to specify the mechanisms that generate the behavioral profile of development on the balance-scale task. The models are disparate, ranging from connectionist implementations (Dawson & Zimmerman, 2003; McClelland, 1989; Shultz, Mareshal & Schmidt, 1994) to production systems (van Rijn, Someren & van der Maas, 2003) to decision trees (Schmidt & Ling, 1996). Typically, these models have attempted to capture the sequence of Siegler's four core rules, and have been judged on their ability to capture the complete range of behavioral phenomena (van Rijn et al., 2003). However, Dawson and Zimmerman (2003) have argued that computational modeling has been preoccupied with fitting the data. Since none of the models give a perfect fit and the detailed data are themselves contested, at this stage the contribution of models should be a qualitative understanding of the mechanisms underlying rule transitions.

Despite the wealth of research on the balance scale task, one area has remained relatively under explored until recently. This is the question of variability. The study of variability in cognitive development is important for three reasons. First, within a single individual, it has been argued that increased variability in performance presages the onset of developmental transitions (Jansen & van der Maas,

2002). Second, variability across individuals of the same age gives a window onto general or specific intelligence. Third, variations in development from the normal pathway are found in disorders, sometimes exhibiting delay, sometimes failure to reach more complex levels of reasoning, and sometimes qualitatively atypical patterns. Implemented models have generally focused on the normative (average) pathway, yet each type of variability must ultimately be explained at a mechanistic level (Thomas & Karmiloff-Smith, 2003).

The following sections report an initial set of simulation results investigating sources of variability in the balance scale task. First we introduce our normal model of development. Second, we explore how changes to the model's computational parameters, representations, and training environment alter its behavioral profile. Third, we evaluate variability in a single case study.

### The Normal Model

The normal model was defined as a 3-layer feedforward connectionist network consisting of an input layer of 20 units representing the number of weights placed (up to 5) on each side of the scale (5 distances either side), a hidden layer of 4 units, and an output layer of 2 units (tip left, tip right). The model used McClelland's (1989) input encoding, where weight and distance information were represented on different units. McClelland's original model separated channels for weight and distance processing channels (i.e., a split hidden layer), a design assumption intended to amplify the model's difficulty in integrating these dimensions. In contrast, we used an undifferentiated network because we wished to avoid using a proprietary network architecture for this particular reasoning problem. There are limitations in our simple model but it remains a useful launching pad to begin an exploration of developmental variability.

The model was trained using back-propagation for 100 epochs, with a learning rate of 0.01. Ten network runs were conducted per manipulation, with initial weights randomized between  $\pm 0.5$ . The standard deviation across runs is depicted in all figures. The training set contained 621 of the possible 625 balance scale problems for a five-peg scale using up to five weights, and was similar to that of McClelland in that balance and weight problems were repeated in the training set. This is based on the original assumption that when learning about balance scale problems, children encounter more experiences that vary on the weight dimension than on the distance dimension. This resulted in a training set consisting of 1069 patterns rather than McClelland's 1125. The remaining 24 problems were used to assess novel performance, which was assessed at 10, 20, 25, 30, 35, 40, 50, 60, 70, 80, 90, and 100 epochs.

The test set consisted of 4 problems from the 6 problem types (see Figure 1). The model's performance on the test set was assessed with 7 test metrics. The metrics captured behavior in line with the following rules: (i) *Rule I*, (ii) *SDD* rule, (iii) *Rule II*, (iv) *QP* rule, (v) *Rule III*, (vi) *addition* rule, and (vii) *Rule IV*. Each metric calculated the percentage of responses consistent with its rule. Note that a

given correct response may be consistent with several rules. For example, Figure 2 shows the problem-space and the proportion of patterns consistent with the four core rules.

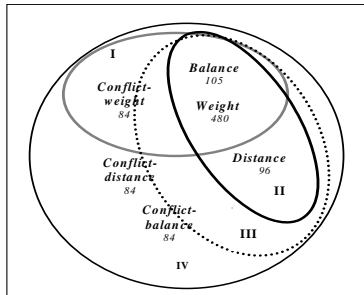


Figure 2: The problem space for rules I to IV

The normal network learned the training set to an accuracy of 98.0% (SD 0.0%). The mean performance of the normal model on each of the test metrics across training is shown in Figure 3 (1HL). Given that we did not separate distance and weight information in the architecture, the network did not exhibit strong evidence of early *Rule I*, *SDD*, or *Rule II* behavior, confirming that weight-distance integration difficulties require architectural assumptions. However, our focus here is upon the model’s balance scale behavior around *Rule III*, since much of the literature has focused on this phase. The sequence of metrics that best characterized the development of the model was: *QP* -> *Rule III* -> *addition rule* -> *Rule IV*.

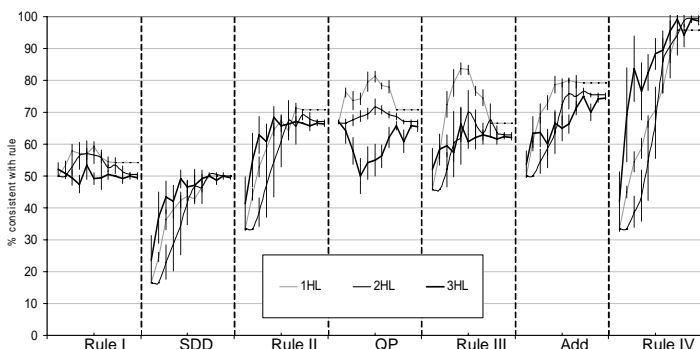


Figure 3: Developmental phases of the normal model ‘1HL’ and networks with 2 and 3 Hidden Layers (2HL, 3HL)

### Exploring Variability

Variability was explored by making a series of systematic changes either to the normal model’s computational parameters, to the problem encoding, or to its environment.

#### Variability and Computational Parameters

We varied (i) the number of hidden layers, (ii) the number of hidden units in a single layer, and (iii) the learning rate.

**Increasing the number of hidden layers** The performance of the model was tested with 2 and 3 hidden layers (HL), with 4 units per layer. Additional hidden layers tend to increase the computational complexity of the mappings that can be learned by a network, while slowing down learning since the error signal must filter back through more levels.

So that learning would fall within a 100-epoch window, the learning rate (*lr*) was increased as follows: 1HL=0.01, 2HL=0.02, 3HL=0.2 (these values hold for subsequent use of these architectures unless otherwise stated)<sup>1</sup>. These networks achieved mean accuracy levels on the training set of 98.0, 99.8, and 100.0% respectively. The developmental profiles of the networks are included in Figure 3. Increasing the number of hidden layers altered the number of transitions in behavior made prior to approaching *Rule IV* performance (we define a transition as a shift in the rule that covers the most behavior). The standard deviation across runs went up as the number of hidden layers increased but notably, phases of development became less incremental. For example, the sequence of closest fitting metrics for models with 2HL was: *QP* -> *addition rule* -> *Rule III* -> *Rule IV*, but was just *QP* -> *Rule IV* for networks with 3HL. (This pattern did not result from *lr* changes, since it did not arise when 1HL was trained with a learning rate of 0.2). Increasing the power of the network reduced the number of transitional states it went through in reaching mastery.

#### Increasing the number of hidden units in a single layer

Expanding the number of units in a single layer increases the capacity of the network to learn more patterns of a given complexity, and allows it to learn a given problem with smaller weights, thereby requiring less learning. We evaluated networks with 4, 10, and 20 units in the hidden layer for the normal 1HL network. After training, the all networks had a mean accuracy of 98.0%. Their developmental profiles are shown in Figure 4. Increasing the number of hidden units did not change the profiles compared to the normal case. We explored this manipulation in the 2HL and 3HL networks and found the same result. If the capacity of the system is measured in parallel processing resources, additional capacity did not alter the transitional stages through which the system passed but altered the rate at which it did so.

#### Reducing the learning rate

Individual differences and developmental disorders are sometimes characterized in terms of *delay*. This term is usually descriptive, but one obvious way to implement it is to turn down the learning rate. This would not explain why delay is frequently uneven across problem domains, but we can at least address how learning rate alters the transitions that the system exhibits. Learning rate was reduced in the normal network in four steps as follows: 0.08, 0.06, 0.04, and 0.02. After 100 epochs, these networks achieved mean accuracies 98.0, 97.1, 94.1, and 56.7% respectively. Figure 5 depicts their developmental phases, with the four steps labeled LR1 to LR4. Slower learning rates caused roughly parallel shifts for all metrics from right to left. That is, while development slowed down, the order of the transitions between types of reasoning behavior remained the same. When the learning

<sup>1</sup> These networks showed qualitatively equivalent results to multi-hidden-layer networks trained on the same *lr* with an extended training time.

rate was insufficient to achieve mastery within the fixed time window of 100 epochs, performance terminated at a less complex level (e.g., LR3 terminated at *Rule III* rather than *IV*, LR4 at the *addition* rule). However, were training extended, *Rule IV* would be reached in both cases. By contrast, developmental disorders typically exhibit asymptoting performance at less complex levels of reasoning. For individual differences, it is unclear whether everyone eventually ‘catches up’. Reduced learning rate does not, therefore, seem a good (sole) candidate to explain the type of developmental delay found in disorders.

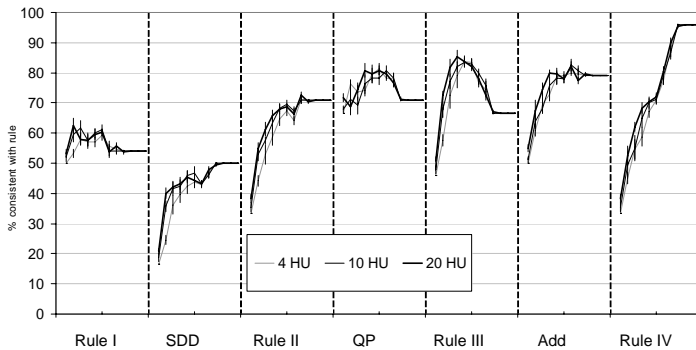


Figure 4: Profile for models with 4 (normal), 10 and 20 hidden units in a single hidden layer

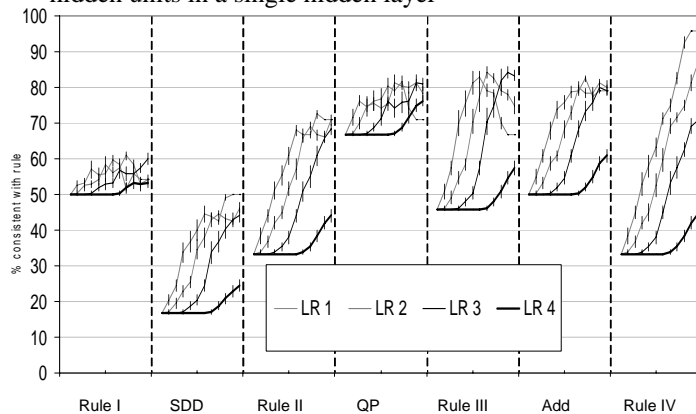


Figure 5: The 1HL model with reducing learning rates

### Variability and the Problem Encoding

We explored two variations in the problem encoding. These correspond to alterations in the way in which the problem is presented to the child (perhaps in the salience of different information or options) or to alterations in how the problem is encoded in the part of the cognitive system required to predict outcomes of balance-scale problems. We either: (i) added a further response option so that the scale could either tip left, tip right, or balance; or (ii) altered the input coding so that information about the weights was represented with position-specific units.

**Changing the response options** In the normal model, there are two output units whose activation can vary between 0 and 1. If the left output unit is more active than the right unit by more than 0.33, the response is ‘tip left’, and vice versa for ‘tip right’. If the difference between the units is less than

0.33, the response is taken to be ‘balance’ (McClelland, 1989). However, since balance is a legitimate response for a proportion of the problems, it could reasonably be encoded as a separate output unit. For this condition, a response was considered correct if the activation of the corresponding output unit was  $\geq 0.5$  and the activation of any other output unit was  $< 0.5$ . Finally, because encoding of the problem domain could alter its complexity, we contrasted performance on 1HL, 2HL, and 3HL networks with 4 hidden units per layer. After training, these networks achieved mean accuracy levels of 86.3, 87.5, and 99.2% respectively. The developmental phases are shown in Figure 6. Comparison with Figure 3 reveals that the additional response option dramatically changed the pattern of transitions. 1HL and 2HL networks began in *Rule I* and did not exceed *Rule II*. Only the 3HL network reached *Rule IV*. Changing the response options altered the categorization that the internal representations had to make across problem types. It ramped up the complexity of the task since balances must be computed internally rather than left to the competition between left and right output units.

**Combined Weight-Distance Encoding** In McClelland’s (1989) formulation, weight and distance information were encoded separately. However, one could represent the amount of weights on each peg locally at each distance. For this manipulation, there were 10 input units, one for each peg on the balance scale. The activation level coded the number of weights placed on a peg. Activation ranged from 0 to 1 and each weight was represented by an increment of 0.2. Thus, three weights on a peg corresponded to an activation of 0.6. The composition of the training set and the output responses remained as normal. Networks with 1, 2, and 3 hidden layers were run to assess the demands of this encoding. The results are in Figure 7.

The final performance of the models was poorer than with the normal encoding by around 20% (1HL=80.2%, 2HL=65%, 3HL=90.3%). As above, only the 3HL network achieved *Rule IV* reasoning as the closest fitting metric at the end of training. 1HL only reached the *QP* rule. Again, a recoding of the problem domain, this time at input, increased the complexity of the task and altered the developmental phases exhibited by the model.

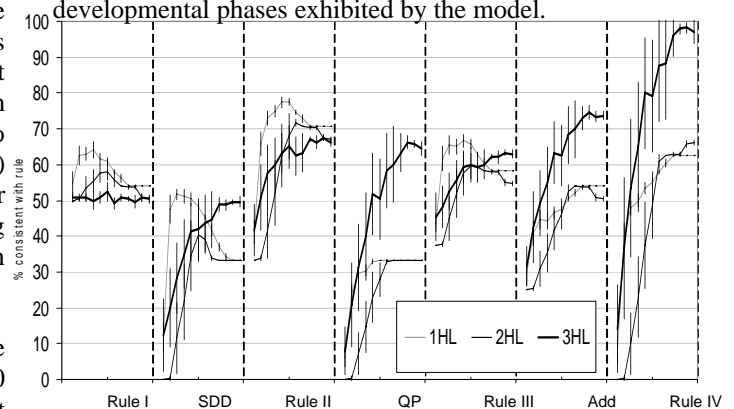


Figure 6. Profile over training for models with 3 response options, shown for 1HL, 2HL, and 3HL networks

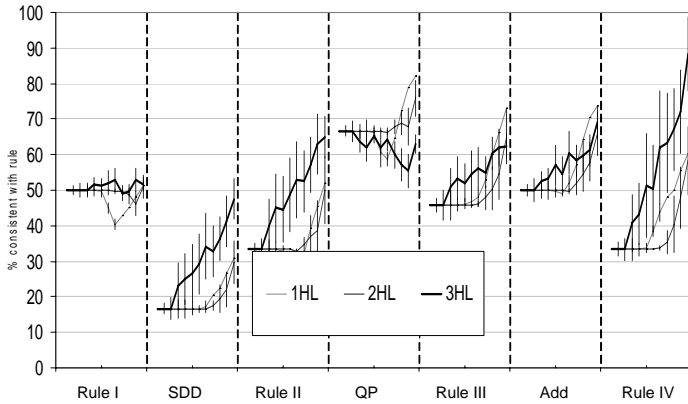


Figure 7: Profile over training for models trained on an environment with combined encoding

### Variability and the Engaged Environment

Since development in the balance scale task corresponds to the child’s active exploration of the domain, we refer to the training set as the *engaged environment*. We created two variations in the environment: (i) a training set without the weight dimension bias, and (ii) an impoverished training set with restricted coverage of the problem space. In these cases, the normal architecture and problem encoding was used. 1HL, 2HL, and 3HL networks with 4 units per layer were also contrasted to explore whether additional representational power could overcome limitations in the engaged environment.

**An Impoverished Engaged Environment** This engaged environment consisted of a subset of 703 training patterns, which excluded any problems where the distances from the fulcrum on *both* sides were  $\geq 3$ . After training, the 1HL, 2HL, and 3HL reached accuracy levels of 97.8, 99.7, and 100.0% respectively. This environment had an adverse effect on the single hidden layer network, where the closest fitting metric at the end of training was *Rule III* (Figure 8) instead of the normal *Rule IV*. The number of closest fitting metrics was also fewer across training, indicating fewer transitions. In contrast, for 2HL and 3HL networks, the closest fitting test metric at the end of training was *Rule IV*, with the 2HL network making more transitions than the 3HL. For all models, there was a considerable increase in variability between individual runs compared to the normal environment. This impoverished environment, then, increased developmental variability between individuals but, importantly, could be compensated for in a more powerful learning system with respect to this test set.

**An Unbiased Engaged Environment** The unbiased engaged environment consisted of 1069 patterns where the original bias for the weight dimension was removed. The duplicated weight problems were replaced with a random selection of patterns already in the training set. All models trained using this environment were able to reach *Rule IV* performance. However, this environment reduced the number of transitions between rules across training. The developmental phases for 1HL, 2HL, and 3HL networks are depicted in Figure 9.

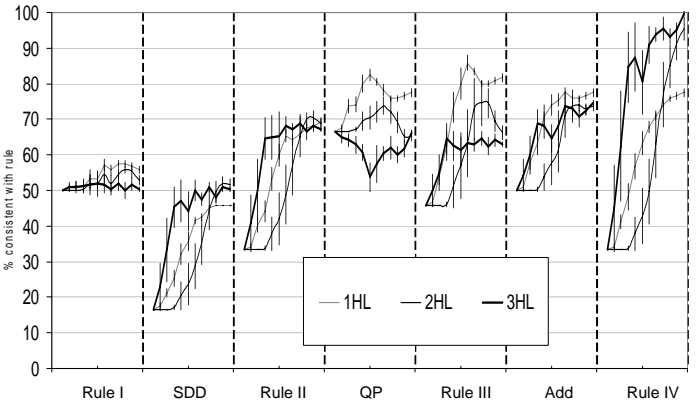


Figure 8: Developmental phases for models in an impoverished engaged environment

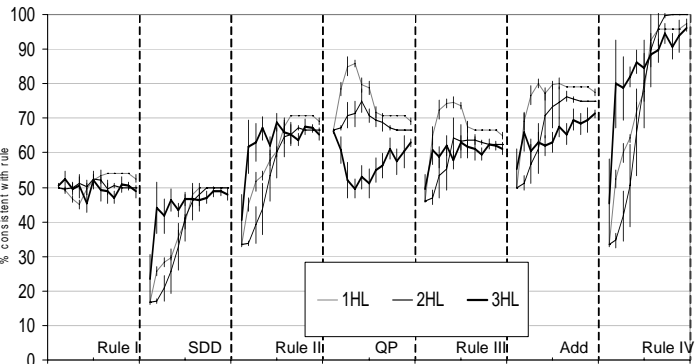
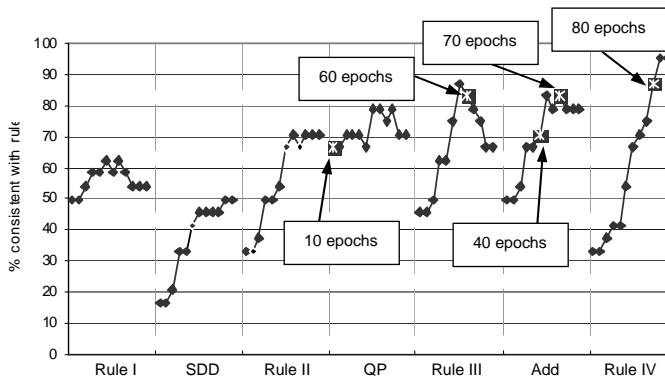
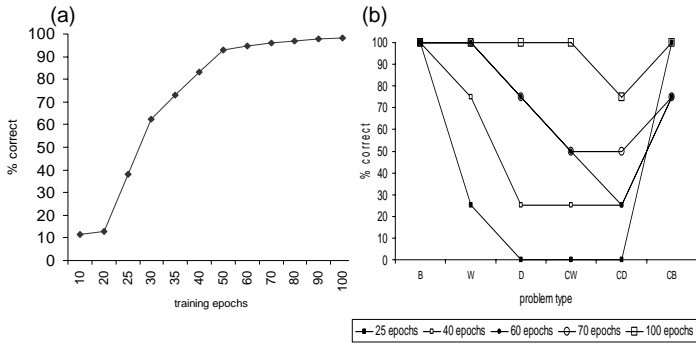


Figure 9: Profile over training for models trained on an unbiased engaged environment.

### Individual Variability: A Case Study

Variability also occurs during the development of individual children, including regression to less sophisticated rules. However, averaging across individuals risks producing variability not found in any one, which may be the case for simulations as well. In this section, we report the rule transitions in the trajectory of a single network (1HL,  $lr=0.008$ , normal encoding and training set). Performance on the training set is shown in Figure 10(a), while 10(b) illustrates performance on the 6 problem types in the test set at 25, 40, 60, 70, and 100 epochs. Figure 11 depicts the rule transitions shown by this individual network. The model made the following transitions: *QP* -> *addition* -> *Rule III* -> *addition* -> *Rule IV*. The trajectory confirms variability around *Rule III*, with a jump from *QP* to *addition*, back to the less sophisticated *Rule III*, returning to *addition*, and on to *Rule IV*. Inspection of Figure 10(b) suggests that this variability is driven by the network’s attempts to solve the low-salience distance problems. Balance and weight problems are performed well from early on, but the network struggles to accommodate distance and conflict-distance problems, inducing greater variability and more transitions between 60 and 80 epochs. In sum, the variability found in averaged data is not an artifact of averaging but found in individual runs. Apparent rule transitions, including regressions, are a key part of the network’s attempts to integrate weight and distance information in solving balance scale problems.



## Conclusions

Mechanisms underlying variability in cognitive development are important for understanding individual differences and developmental disorders, as well as normative development. Simulation of the balance scale task indicated that variations of internal computational parameters, problem encoding, and engaged environment all act on the complexity of the reasoning exhibited by the network during learning, including the findings that more hidden layers increase complexity but not more units per layer (contrasting reasoning power with capacity); that a slower learning rate does not reduce complexity per se, and is therefore a poor model of unresolved developmental delay; and that an impoverished environment can reduce complexity but be compensated for by a more powerful learning system.

## References

Dawson, M. R. W., & Zimmerman, C. (2003). Interpreting the internal structure of a connectionist model of the balance scale task. *Brain and Mind, 4*, 129-149.

Ferretti, R. P., & Butterfield, E. C. (1986). Are children's rule-assessment classifications invariant across instances of problem types? *Child Development, 57*, 1419-1428.

Halford, G. S., Andrews, G., Dalton, C., Boag, C., & Zielinski, T. (2002). Young children's performance on the balance scale: The influence of relational complexity. *Journal of Experimental Child Psychology, 81*, 417-445.

Hardiman, P. T., Pollatsek, A., & Well, A. D. (1986). Learning to understand the balance beam. *Cognition & Instruction, 3(1)*, 63-86.

Inhelder, B., & Piaget, J. (1958). *The growth of logical thinking from childhood to adolescence*. NY: Basic Books.

Jansen, B., & van der Maas, H. (2002). The development of children's rule use on the balance scale task. *Journal of Experimental Child Psychology, 81*, 383-416.

Jansen, B., & van der Maas, H. (1997). Statistical test of the rule assessment methodology by latent class analysis. *Developmental Review, 17*, 321-357.

McClelland, J. L. (1989). Parallel distributed processing: Implications for cognition and development. In M. G. M. Morris (Ed.), *Parallel distributed processing, implications for psychology and neurobiology* (pp. 8-45). Oxford: Clarendon Press.

Normandeau, S., Larivée, S., Roulin, J. L., & Longeot, F. (1989). Young children's knowledge of balance scale problems. *Journal of Genetic Psychology, 148*, 79-94.

Schmidt, W. C., & Ling, C. X. (1996). A decision-tree model of balance scale development. *Machine Learning, 24*, 203-230.

Shultz, T. R., Mareschal, D., & Schmidt, W.C. (1994). Modeling cognitive development on balance scale phenomena. *Machine Learning, 16*, 57-86.

Siegler, R. S. (1976). Three aspects of cognitive development. *Cognitive Psychology, 8*, 481-520.

Siegler, R. S. (1981). Developmental sequences within and between concepts. *Society for Research in Child Development Monographs, 46* (Whole no. 189).

Siegler, R. S. (2002). Microgenetic studies of self-explanations. In N. Granott & J. Parziale (Eds.), *Microdevelopment: Transition processes in development and learning* (pp. 31-58). NY: Cambridge University.

Siegler, R. S., & Chen, Z. (2002). Development of rules and strategies balancing the old and the new. *Journal of Experimental Child Psychology, 81*, 446-457.

Thomas, M. S. C. & Karmiloff-Smith, A. (2003). Connectionist models of development, developmental disorders and individual differences. In R. J. Sternberg, J. Lautrey, & T. Lubart (Eds.), *Models of intelligence: International perspectives*, (p. 133-150). APA.

van Maanen, L., Been, P., & Sijtsma, K. (1989). The linear logistic test model and heterogeneity of cognitive strategies. In E. E. Roskam (Ed), *Mathematical psychology in progress* (pp. 267-288). Berlin: Springer.

van der Maas, H., & Jansen, B. (2003). What response times tell of children's behavior on the balance scale task. *Journal of Experimental Child Psychology, 85*, 141-177.

van Rijn, H., van Someren, M., & van der Maas, H. (2003). Modeling developmental transitions on the balance scale task. *Cognitive Science, 27*, 227-257.

Wilkening, F., & Anderson, N. H. (1982). Comparison of two rule-assessment methodologies for studying cognitive development and knowledge structure. *Psychological Bulletin, 92*, 215-237.

**Acknowledgements:** This research was supported by UK MRC CE Grant G0300188 to Michael Thomas