

## А Н Н О Т А Ц И Я

В предыдущих работах была предложена модель работы мозга, формализующая и объясняющая Теорию Функциональных Систем П.К.Анохина. Основываясь на этих работах в дальнейшем была предложена схема работы функциональных систем и объяснение Теории Движений Н.А.Бернштейна. В результате была разработана некоторая модель работы мозга. Настало время экспериментальной проверки работоспособности теории.

Цель данной работы в том, что бы, основываясь на полученной модели, построить некоторый простейший анимат и показать, что он обучается достаточно эффективно и, по крайней мере, не хуже, чем на основании существующих подходов, основанных на нейронных сетях и потактовом обучении (Reinforcement Learning).

Такой анимат создан, запрограммирован и проведены машинные эксперименты по его обучению. Результаты сравнения его поведения, основанные на предложенной модели и моделях, основанных на Reinforcement Learning, показали, что предложенная модель обучается и действует эффективней.

## СИСТЕМА УПРАВЛЕНИЯ АНИМАТОМ ОСНОВАННАЯ НА ТЕОРИИ ФУНКЦИОНАЛЬНЫХ СИСТЕМ П.К.АНОХИНА<sup>1</sup>

А. В. ДЕМИН<sup>1</sup>, Е. Е. ВИТЯЕВ<sup>2</sup>

<sup>1</sup>Институт систем информатики имени А. П. Ершова СО РАН г. Новосибирск

<sup>2</sup>Институт математики СО РАН г. Новосибирск, e-mail: [vityaev@math.nsc.ru](mailto:vityaev@math.nsc.ru)

### 1. Введение.

В работах [1-6] была предложена модель работы мозга, формализующая и объясняющая Теорию Функциональных Систем (ТФС) П.К.Анохина. Основываясь на этих работах в дальнейшем была предложена схема работы функциональных систем [7] и объяснение Теории Движений Н.А.Бернштейна [8]. В результате была разработана некоторая модель работы мозга. Настало время экспериментальной проверки работоспособности теории.

Цель данной работы состоит в том, что бы, основываясь на работах [1-9], построить некоторый простейший анимат и показать, что он обучается достаточно эффективно и, по крайней мере, не хуже, чем на основании существующих подходов, основанных на нейронных сетях и потактовом обучении (Reinforcement Learning).

### 2. Система управления аниматом.

В соответствии с Теорией Функциональных Систем будем считать, что моделируемая система управления аниматом имеет иерархическую архитектуру, в которой базовым элементом системы управления является отдельная функциональная система (ФС). При такой архитектуре, ФС верхнего уровня ставят цели системам нижнего уровня. При этом можно считать, что каждая ФС решает задачу достижения цели, используя те же методы, что и остальные ФС. На рис.1 приведена архитектура системы управления аниматом.

Задача отдельной функциональной системы:

- При заданной цели (подцели) и известной информации об окружающей среде и состоянии ФС, найти оптимальный способ достижения цели.

---

<sup>1</sup> Эта работа поддержана грантом РФФИ 05-07-90185в, Интеграционным проектом СО РАН №119, Программой президента Российской Федерации поддержки научных школ 2112.2003.1

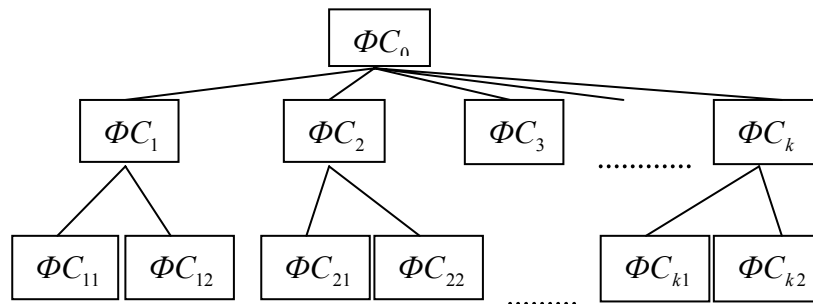


Рис. 1. Архитектура системы управления

- Если на основе прогноза найдено действие, обеспечивающее достижение цели, то дать команду на выполнение этого действия.
- Осуществить контроль правильности выполнения действия, т.е. проверить соответствие параметров достигнутого и желаемого результатов.

### 3. Модель работы функциональной системы.

На рис.2 приведена модель работы функциональной системы, основанная на работе [1]. Будем считать, что в некоторый момент времени ФС ставится цель  $P_0$ . Цель ставится в виде запроса к ФС – достичь цель  $P_0$ . На вход ФС также подается информация об окружающей среде в виде описания ситуации  $P_{i1}, \dots, P_{im}$ . В процессе афферентного синтеза из памяти извлекается вся информация, связанная с достижением цели  $P_0$ . Эта информация хранится в памяти в виде множества закономерностей вида  $\langle P_{i1}, \dots, P_{ik}, A_i \rangle \rightarrow P_0$ . Поскольку

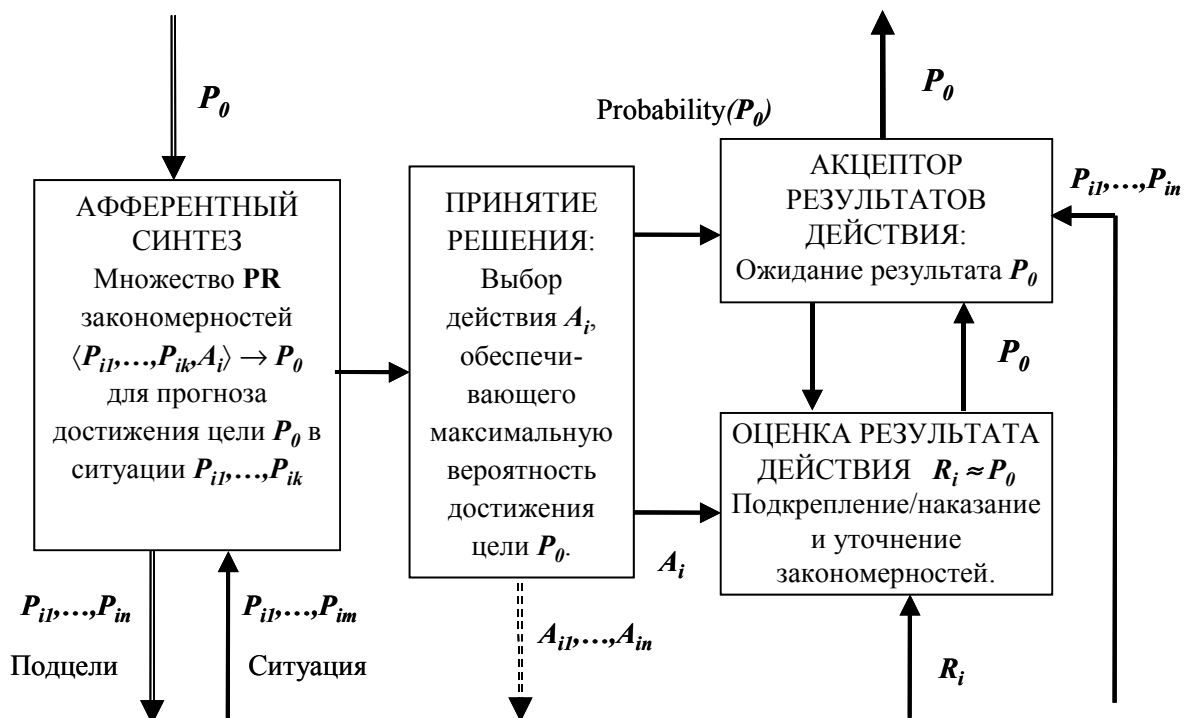


Рис 2. Модель работы функциональной системы

информация об окружающей среде уже задана в виде описания ситуации  $P_{i1}, \dots, P_{im}$  еще до постановки цели, то из памяти автоматически извлекается только та информация, связанная с достижением цели, которая может быть применена в данной ситуации. Это достигается использованием (извлечением из памяти) только тех закономерностей, в которых свойства ситуации выполнены, то есть, свойства ситуации  $P_{ij}$ , которые есть в условии  $P_{i1}, \dots, P_{ik}$  закономерности должны содержаться в описании ситуации  $P_{i1}, \dots, P_{im}$ .

Среди условий  $P_{i1}, \dots, P_{ik}$  закономерности содержатся не только свойства ситуации, но и подрезультаты  $P_{i1}, \dots, P_{in}$ , достижение которых необходимо для достижения цели. Достижение подцелей осуществляется отправкой запроса на их достижение вниз по иерархии, что обозначено на рисунке двойной стрелкой вниз. Эти запросы активируют в ФС более низкого уровня всю информацию, связанную с достижением этих подцелей. Достижение подцелей может потребовать достижение других подцелей в иерархии целей и т.д. Если какая-то из подцелей не может быть выполнена в данной ситуации (нет закономерностей предсказывающих достижение подцели в данной ситуации), то ответом на запрос будет отказ и соответствующая закономерность исключается из рассмотрения.

Активация закономерностей  $\langle P_{i1}, \dots, P_{ik}, A_i \rangle \rightarrow P_0$  в блоке афферентного синтеза автоматически извлекает из памяти тот набор действий  $A_i$ , включая действия, требуемые для достижения подцелей, которые могут привести к достижению цели  $P_0$ . Весь этот набор действий вместе с оценками условных вероятностей достижения цели и подцелей передается в блок принятия решений. Блок принятия решений просматривает все действия  $A_i$  вместе с активирующими их закономерностями  $\langle P_{i1}, \dots, P_{ik}, A_i \rangle \rightarrow P_0$  и иерархией подцелей и соответствующих действий, и выбирает такое действие, которое с учетом вероятностей выполнения подцелей дает максимальную оценку вероятности достижения цели. Далее, действие  $A_i$  и все действия, необходимые для достижения подцелей, запускаются на выполнение. В начальной стадии обучения, когда еще нет правил, либо нет ни одного правила, применимого в данной ситуации, действие соответствующей ФС выбирается случайным образом из имеющегося арсенала действий и прогноз отсутствует.

Прогноз ожидаемого результата  $P_0$  и всех подрезультатов для всех подцелей отправляется в акцептор результатов действий. Кроме того, во всех ФС более нижнего уровня прогноз подрезультатов также отправляется в акцептор результатов действия соответствующих подсистем.

Данные о полученном результате  $R_i$  поступают в акцептор результатов действий блока оценки результата. Проводится сравнение спрогнозированного и полученного результа-

тов. В случае совпадения прогноза и результата, с заданной степенью точностью, закономерность, выбранная в блоке принятия решений, подкрепляется, в противном случае, наказывается. Закрепление/наказание состоит в увеличении/уменьшении условной вероятности закономерности. Кроме того, после каждого действия производится уточнение набора правил, как описано в [1]. Если после уточнения для данного состояния находится закономерность с условной вероятностью, большей чем у закономерности, использованной ранее, то новая закономерность будет в дальнейшем использоваться для прогноза и принятия решения.

Семантический вероятностный вывод [10] позволяет найти набор  $PR$  закономерностей вида  $\langle P_{i1} \& \dots \& P_{ik} \& A_i \rangle \rightarrow P_0$ , с максимальной условной вероятностью предсказывающий результат  $P_0$  действия  $A_i$  в состоянии  $\langle P_1, \dots, P_k \rangle$ .

#### 4. Иерархия функциональных систем

Представим функциональные системы более схематично рис.3, и в качестве иерархии ФС приведем два уровня ФС. Функциональные системы не являются раз и навсегда за-

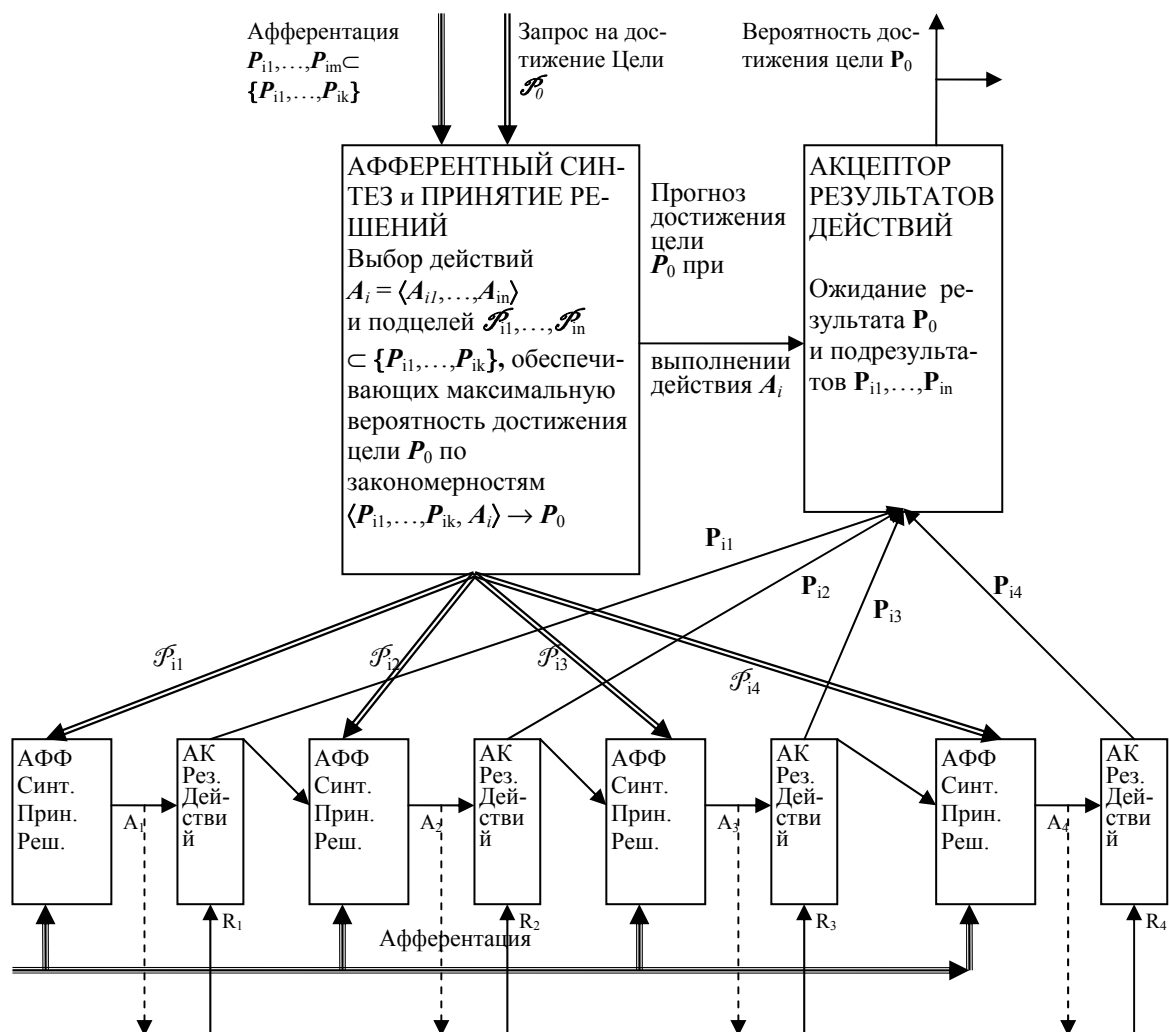


Рис.3 Два уровня иерархии функциональных систем

данными образованиями. Они меняются и формируются в зависимости от целей. Цели и подцели в свою очередь тоже формируются в зависимости от успешности достижения конечных целей. Покажем как с помощью закономерностей могут автоматически формироваться цели и подцели.

Расширим понятие результата так, чтобы он мог автоматически формироваться в процессе целенаправленных действий в сложной вероятностной среде:

- а) результат должен обладать свойством ветвления: если получен некоторый результат, то дальнейшие действия могут определяться неоднозначно;
- б) результат должен содержать набор признаков, которые определяют, что цель цепочки действий достигнута и можно переходить к одной из следующих цепочек действий, т.е. результат - это фиксация законченности действия, обеспечивающая возможность осуществления некоторого следующего действия.

Условие а) определения результата естественным образом улавливается закономерностями, так как закономерности хорошо прогнозируют результат последовательности некоторых элементарных действий (данного уровня), если эта последовательность действий «стандартна» (начавшись, она продолжается до некоторого результата без изменений). В этом случае с большой вероятностью закономерности прогнозируют выполнение цепочки действий до получения результата. На рис.3 это действия  $A_1, A_2, A_3, A_4$ , приводящие к результатам  $R_1, R_2, R_3, R_4$ . Акцептор результатов действия сличает результаты  $R_1, R_2, R_3, R_4$  с предсказанными по закономерностям и, в случае совпадения, выдает ответы  $P_1, P_2, P_3, P_4$  на запросы  $\mathcal{R}_1, \mathcal{R}_2, \mathcal{R}_3, \mathcal{R}_4$ . Ответы о достижении цели передаются на входы других блоков (стрелки от акцепторов результатов действий вниз к следующим блокам). Эти ответы автоматически будут включаться в условия закономерностей последующих действия, так как сигнал о том, что предыдущее действие завершено, увеличивает вероятность завершения последующего действия. По определению семантического вероятностного вывода [10] любой сигнал, увеличивающий вероятность прогноза, автоматически включается в условие закономерности.

Условие б) также выполнено, так как сигналы от обратной афферентации, свидетельствующие о действительном завершении предыдущего действия, увеличивают вероятность достижения результата следующего действия. Более детально, на нейронном уровне, процесс автоматического формирования целей и подцелей описан в [4].

## 4. ОПИСАНИЕ МОДЕЛИ

**1. Описание модели.** Приведем схему работы анимата, реализующую схему функциональных систем рис.3 [8]. Будем предполагать, что система управления аниматом функ-

ционирует в дискретном времени  $t = 0, 1, 2, \dots$ . Пусть анимат имеет некоторый набор сенсоров  $S_1, S_2, \dots, S_n$ , характеризующих состояние внешней и внутренней среды, и набор возможных действий  $A_1, A_2, \dots, A_m$ . Среди множества сенсоров выделим сенсор  $SA$ , который представляет информацию о совершенном действии. Считаем, что история деятельности анимата хранится в таблице данных  $X = \{X_1, \dots, X_t\}$ , в которой  $t$ -я строка таблицы содержит показания сенсоров в момент времени  $t$ :  $X_t = \{S_1^t, S_2^t, \dots, S_n^t, SA_t\}$ ,  $S_1^t, S_2^t, \dots, S_n^t$  - значения сенсоров  $S_1, S_2, \dots, S_n$  в момент времени  $t$ . На множестве  $X$  определим множество предикатов  $P = \{P_1(t), \dots, P_k(t), PA_1(t), \dots, PA_m(t)\}$ , где  $P_i(t)$  - *сенсорные предикаты*, определяющие некоторые условия на показания сенсоров в момент времени  $t$ ;  $PA_i(t) \Leftrightarrow (SA(t) = A_i)$  - *активирующие предикаты*, показывающие, что в момент времени  $t$  было совершено действие  $A_i$ .

Введем понятие *предиката-цели*  $PG(t) = P_{i1}(t) \& P_{i2}(t) \& \dots \& P_{in}(t)$ , реализующего условие достижения цели в момент времени  $t$ .

Каждой функциональной системе  $\Phi C_j$  соответствует некоторая цель  $G_j$ , достижение которой является задачей данной ФС, и предикат-цель  $PG_j$ , характеризующий условие достижения цели.

Каждая  $\Phi C_j$  содержит свой набор предикатов  $P_j = P \cup \{PG_{j1}, \dots, PG_{jn}\}$ , где  $PG_{jk}$  - предикаты-цели, соответствующие целям нижестоящих по иерархии ФС, подчиненных данной ФС. Каждая  $\Phi C_j$  содержит множество закономерностей  $PR_j$  вида:  $\{P_{i1}, \dots, P_{ik}, PA_i \mid PG_{j1}, \dots, PG_{jn}\} \rightarrow PG_j$ , где  $\mid$  - означает, что в условии правила стоит только одно из указанных в фигурных скобках выражений. Каждая такая закономерность характеризуется некоторой оценкой  $p$  вероятности достижения цели  $PG_j$ , при выполнении условия закономерности.

Предположим, что в некоторый момент времени  $t$  система  $\Phi C_j$  получила запрос на достижение цели  $PG_j$ . Тогда из множества закономерностей  $PR_j$  извлекаются все закономерности, условие которых выполнено в текущий момент времени  $t$ . Если условие закономерности содержит предикаты-подцели  $PG_{j1}, \dots, PG_{jn}$ , то ФС отправляет запрос на достижение этих подцелей вниз по иерархии ФС. Среди всех отобранных закономерностей выбирается та закономерность, которая с учетом вероятностей выполнения подцелей дает максимальную оценку  $f$  вероятности достижения цели. Оценка  $f$  закономерности  $\{P_{i1}, \dots, P_{ik}, PA_i \mid PG_{j1}, \dots, PG_{jn}\} \rightarrow PG_j$  вычисляется следующим образом:  $f(PG_j \mid \{PS_{i1}, \dots, PS_{ik}, PA_i \mid PG_{j1}, \dots, PG_{jn}\}) = p \cdot f(PG_{j1}) \cdot \dots \cdot f(PG_{jn})$ , где  $p$  - оценка вероятности данной закономерности,  $f(PG_{jk})$  - оценки вероятностей достижения подцелей. Если все

условия выбранной закономерности выполнены, то действие  $A_i$  запускается на выполнение. Если множество закономерностей  $PR_j$  пусто, либо нет ни одной закономерности, применимой в данной ситуации, то действие выбирается случайно из арсенала имеющихся действий. После совершения действия обновляются показания сенсоров, оценивается результат действия и уточняется набор правил  $PR_j$  (см. ниже).

**2. Оценка результатов действий.** Предположим, что системой  $\Phi C_j$  получен запрос на достижение цели  $G_j$  и после достижения цели был получен результат  $R_j$ . Определим оценку результата действий. Если после очередного действия предикат цели ложен  $PG_j = 0$ , то результат не достигнут и оценка результата действий  $d^j(t) = 0$ . Если результат был получен в момент времени  $t_0$ , то все оценки  $d^j(t)$ , начиная с момента времени  $t_0$  и до предыдущего момента достижения цели  $t_1$ , пересчитываются следующим образом:  $d^j(t) = re^{\alpha(t-t_0)}$ ,  $t_1 < t \leq t_0$ , где  $r$  – функция оценки качества полученного результата,  $\alpha > 0$  – параметр  $r = \begin{cases} 0, & \text{если } PG_j = 0 \\ \|G_j - R_j\|, & \text{если } PG_j = 1 \end{cases}$ , где  $\|\dots\|$  – мера близости между полученным результатом  $R_j$  и поставленной целью  $G_j$ . Каждая  $\Phi C_j$  хранит оценки результатов действий  $d^j(t)$  для каждого момента времени  $t$ .

**3. Генерация правил.** Для получения множества закономерностей  $PR_j$ , которые использует система  $\Phi C_j$ , воспользуемся семантическим вероятностным выводом [10].

Семантический вероятностный вывод позволяет находить все закономерности вида  $P_1, \dots, P_m \rightarrow P_0$  с максимальной вероятностью предсказывающие предикат  $P_0$ . Вывод осуществляется на некотором множестве обучающих данных  $Y$  с использованием заданного множества предикатов  $\{P_1, \dots, P_m\}$ .

Данный метод основывается на следующем определении вероятностной закономерности, предложенном в работе [12].

Правило  $P_1, \dots, P_m \rightarrow P_0$  является *закономерностью*, если оно удовлетворяет следующим условиям:

1.  $p(P_1, \dots, P_m) > 0$ ,
2.  $\forall \{P_{i_1}, \dots, P_{i_k}\} \subset \{P_1, \dots, P_m\} \quad p(P_0 | P_1, \dots, P_m) > p(P_0 | P_{i_1}, \dots, P_{i_k})$ .

Здесь  $p$  – оценка условной вероятности правила.

Введем понятие *уточнения* правила. Правило  $P_{i_1}, \dots, P_{i_m}, P_{i_{m+1}} \rightarrow P_0$  является *уточнением* правила  $P_1, \dots, P_m \rightarrow P_0$ , если оно получено добавлением в посылку правила  $P_1, \dots, P_m \rightarrow P_0$  произвольного предиката  $P_{i_{m+1}}$ , и  $p(P_0 | P_{i_1}, \dots, P_{i_{m+1}}) > p(P_0 | P_1, \dots, P_m)$ .



Алгоритм семантического вероятностного вывода:

- На первом шаге генерируется множество уточнений правила  $\rightarrow P_0$  (т.е. правила с пустой посылкой). Это множество будет состоять из правил единичной длины, имеющих вид  $P_{ij} \rightarrow P_0$ , для которых  $p(P_0 | P_{ij}) > p(P_0)$ .
- На k-ом ( $k > 1$ ) шаге генерируется множество уточнений всех правил, созданных на предыдущем шаге. Т.е. для каждого правила  $P_{i1}, \dots, P_{ik-1} \rightarrow P_0$ , сгенерированного на (k-1)-ом шаге, создается множество правил вида  $P_{i1}, \dots, P_{ik-1}, P_{ik} \rightarrow P_0$ , таких, что  $p(P_0 | P_{i1}, \dots, P_{ik-1}, P_{ik}) > p(P_0 | P_{i1}, \dots, P_{ik-1})$ .
- Проверяется нельзя ли из полученных правил удалить какой-то из предикатов так чтобы при этом условная вероятность правила выросла. Если можно, то такие предикаты удаляются из правила.
- Алгоритм останавливается, когда больше невозможно уточнить ни одно правило,

Для того чтобы избежать генерации статистически незначимых правил, вводится дополнительный критерий – оценка на статистическую значимость. Правила, не удовлетворяющие этому критерию, отсеиваются, даже если они имеют высокую точность на обучающем множестве. Для оценки статистической значимости в алгоритме используется критерий Фишера (точный критерий Фишера для таблиц сопряженности) [13].

Очевидно, что все правила, полученные при помощи данного алгоритма, будут являться закономерностями. На рис.4 представлено дерево вывода, соответствующее описанному процессу.

Чтобы найти все закономерности  $\{P_{i1}, \dots, P_{ik}, PA_i | PG_{j1}, \dots, PG_{jn}\} \rightarrow PG_j$  с максимальной ве-

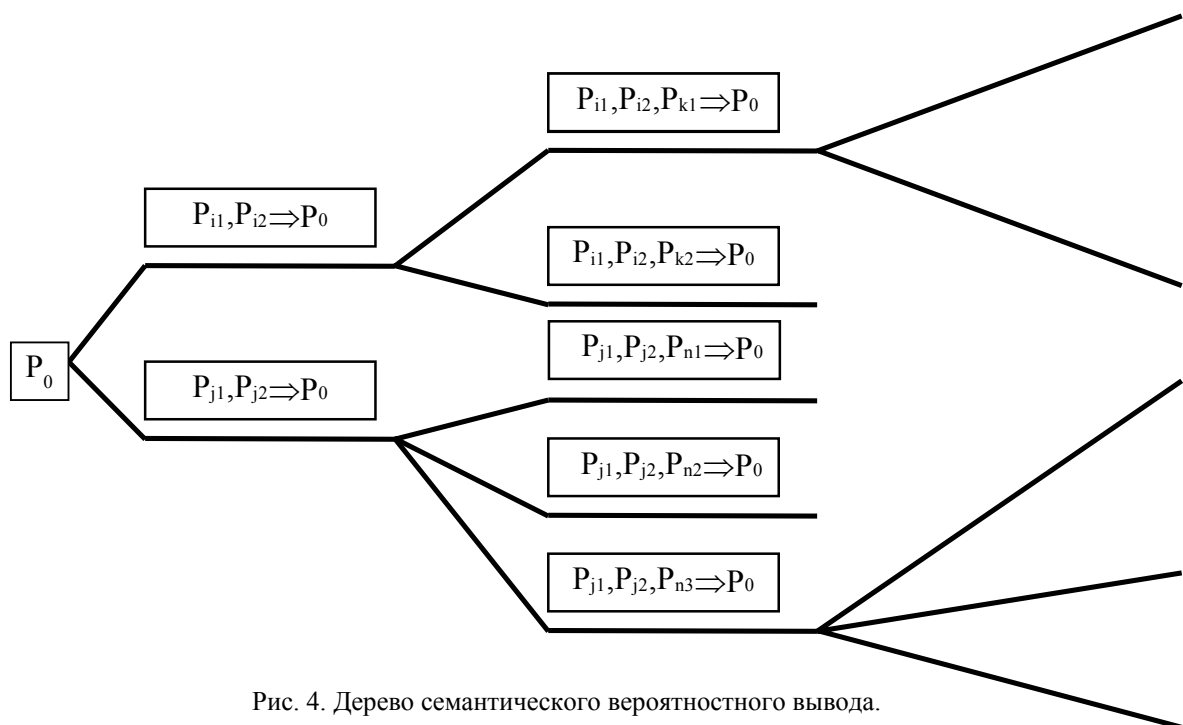


Рис. 4. Дерево семантического вероятностного вывода.

роятностью предсказывающие достижение цели  $G_j$ , строится дерево семантического вероятностного вывода на множестве данных истории деятельности анимата  $X$  и множестве оценок действий  $d^j(t)$  с использованием набора предикатов  $P_j$ , которые использует данная ФС. Оценка условной вероятности  $p$  правила рассчитывается следующим образом:  $p = \sum_{i \in I} d_i^j / \|I\|$ , где  $I$  – множество моментов времени, когда может быть применено данное правило.

**4. Извлечение подцелей.** Изначально система управления аниматом имеет заданную априори иерархию ФС. В простейшем случае она может состоять всего из одной ФС. В процессе деятельности система управления может автоматически выявлять новые подцели и порождать соответствующие ФС. Опишем процедуру порождения новых подцелей и ФС.

Предварительно определим два типа подцелей:

- 1) Подцелями первого типа будем называть ситуации из которых достижение вышестоящей цели прогнозируется одним правилом, содержащим одну цепочку действий, с высокой вероятностью (вероятность становится близкой к 1).
- 2) Подцелями второго типа будем называть ситуации, которые увеличивают вероятность достижения вышестоящей цели, но при этом дальнейшие действия могут определяться неоднозначно.

Для выявления подцелей первого типа среди множества правил  $PR_j$  ФС выбираются правила вида  $R = P_{i1}, \dots, P_{in}, PA_i \rightarrow PG_j$ , имеющие высокую оценку условной вероятности  $p > \delta_1$ , например,  $\delta_1 = 0.9$ . Далее, для каждого отобранного правила  $R$  порождается подцель  $G_i$  и соответствующий предикат-цель  $PG_i = P_{i1} \& P_{i2} \& \dots \& P_{in}$ , равный конъюнкции всех сенсорных предикатов правила  $R$ .

Для выявления подцелей второго типа рассматриваются правила с достаточно высокой оценкой условной вероятности  $p > \delta_2$  (например,  $\delta_2 = 0.7$ ), имеющие вид  $R = P_{i1}, \dots, P_{in}, PA_i \rightarrow PG_j$ . Если среди этих правил найдется хотя бы два правила с разными активизирующими предикатами, но такими, что все сенсорные предикаты одного правила содержатся в другом, то порождается новая подцель  $G_i$  и соответствующий предикат-цель  $PG_i = P_{i1} \& P_{i2} \& \dots \& P_{in}$ , равный конъюнкции всех сенсорных предикатов, содержащихся в обоих правилах.

Таким образом, порождается новая подцель  $G_i$  и соответствующий ей предикат-цель  $PG_i$ , если выполнено одно из следующих условий:

- 1) Если существует правило  $R = P_{i_1}, \dots, P_{i_n}, PA_i \rightarrow PG_j$ , такое, что  $p > \delta_1$ , то формируется подцель  $PG_i = P_{i_1} \& P_{i_2} \& \dots \& P_{i_n}$ .
- 2) Если существуют правила  $R_1 = P_{i_1}, \dots, P_{i_n}, PA_i \rightarrow PG_j$  и  $R_2 = P_{j_1}, \dots, P_{j_m}, PA_j \rightarrow PG_j$ , такие, что  $p_1 > \delta_1$ ,  $p_2 > \delta_2$ ,  $\{P_{i_1}, \dots, P_{i_n}\} \subseteq \{P_{j_1}, \dots, P_{j_m}\}$  и  $A_i \neq A_j$ , то формируется подцель  $PG_i = P_{i_1} \& P_{i_2} \& \dots \& P_{i_n}$ .

Затем, для каждой выявленной подцели  $G_i$  создается новая  $\Phi C_i$ , находящаяся ниже по иерархии и реализующая эту подцель. Для созданной системы  $\Phi C_i$  при помощи семантического вероятностного вывода порождается множество закономерностей  $PR_i$ . Для этого просматривается все множество данных истории анимата  $X$  и выявляются случаи, когда подцель  $G_i$  была реализована, и рассчитывается множество оценок действий  $d^i(t)$   $\Phi C_i$  описанным выше способом. Для всех ФС, находящихся на один уровень выше  $\Phi C_i$ , набор предикатов обогащается еще одним предикатом  $PG_i$ , и генерируются новые правила. Тем самым, множество закономерностей этих ФС обогащаются закономерностями, содержащими новую подцель  $G_i$ .

## 5. ЭКСПЕРИМЕНТ

**1. Описание эксперимента.** Для исследования описанной выше системы управления был поставлен следующий простой эксперимент. При помощи компьютерной программы был смоделирован виртуальный мир и анимат, основной целью которого является обнаружение специальных объектов виртуального мира – «еды». Анимат должен научиться эффективно находить и собирать еду.

Мир анимата представляет собой прямоугольное поле, разбитое на клетки, и содержит три типа объектов: пустые клетки («трава»), препятствия («препятствие»), и еду («еда»). Объекты «препятствие» располагаются только по периметру виртуального мира, образуя тем самым его естественные границы. Анимат может совершать три типа действий: шагнуть на клетку вперед («шаг»), повернуть налево («налево»), повернуть направо («направо»). Когда анимат шагает на клетку, содержащую еду, считается, что он ее «поедает», клетка, на которой находилась еда, очищается и новый объект «еда» случайным образом появляется в другом месте поля. Таким образом, количество еды в виртуальном мире всегда остается постоянным.

Также анимат обладает девятнадцатью сенсорами, девять из которых информируют его о наличии еды («еда на северо-западе», «еда на севере», «еда на северо-востоке», «еда на западе», «еда здесь», «еда на востоке», «еда на юго-западе», «еда на юге», «еда на юго-востоке»), еще девять сенсоров предоставляют информацию о препятствиях («препятст-

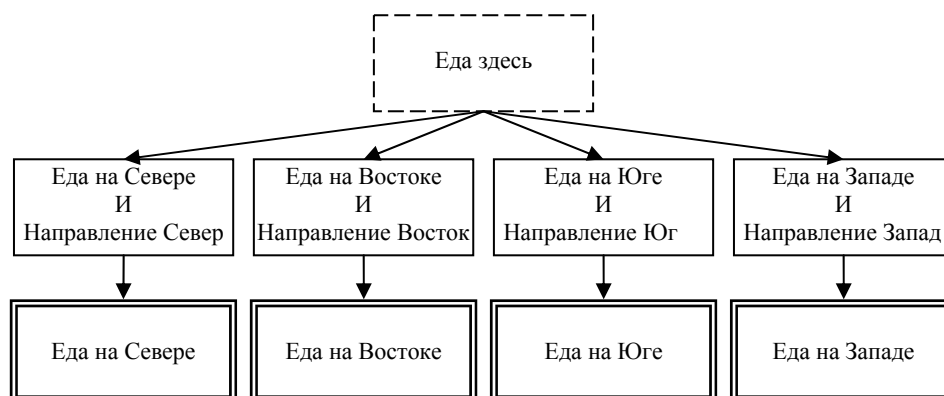


Рис. 5. Иерархия целей, сформированная аниматом в ходе эксперимента

вие на северо-западе», «препятствие на севере», «препятствие на северо-востоке», «препятствие на западе», «препятствие на востоке», «препятствие на юго-западе», «препятствие на юге», «препятствие на юго-востоке») и один сенсор говорит о направлении анимата («направление»). Сенсор направления показывает ориентацию анимата относительно виртуального мира и может принимать следующие значения: «север», «восток», «юг» и «запад». Сенсоры еды и препятствий информируют о наличии данных объектов на клетке, на которой находится анимат, и на соседних с ней клетках, и принимают значения «да» или «нет». Показания этих сенсоров не зависят от ориентации анимата, т.е. сколько бы анимат не крутился на одном месте, их показания не изменяться. Таким образом, что бы эффективно ориентироваться в виртуальном мире, анимат должен научиться сопоставлять свое направление с положением еды и принимать решение о соответствующем действии.

Изначальный набор предикатов анимата состоит из двадцати двух сенсорных предикатов – по одному на каждый сенсор еды и препятствий  $s$ : ( $s = \text{«да»}$ ), и четыре на сенсор направления: («направление» = «север»), («направление» = «восток»), («направление» = «юг») и («направление» = «запад»), и трех активирующих предикатов: ( $A = \text{«шаг»}$ ), ( $A = \text{«налево»}$ ) и ( $A = \text{«направо»}$ ).

Изначально система управления аниматом имеет только одну ФС, целью которой является ощущение наличия еды сенсором «еда здесь», соответствующий предикат-цель имеет вид («еда здесь» = «да»). Когда анимат достигает эту цель, то считается, что он «поедает» еду.

**3. Результаты эксперимента.** Одной из основных задач эксперимента была демонстрация возможности автоматического формирования иерархии целей и результатов при целенаправленном поведении. В ходе эксперимента системой управления аниматом были обнаружены подцели, достижение которых значительно увеличивает вероятность достижения цели. Пример первых двух уровней иерархии целей, сформированной аниматом, представлен на рис.5. На рисунке пунктирной линией обозначена основная цель анимата,

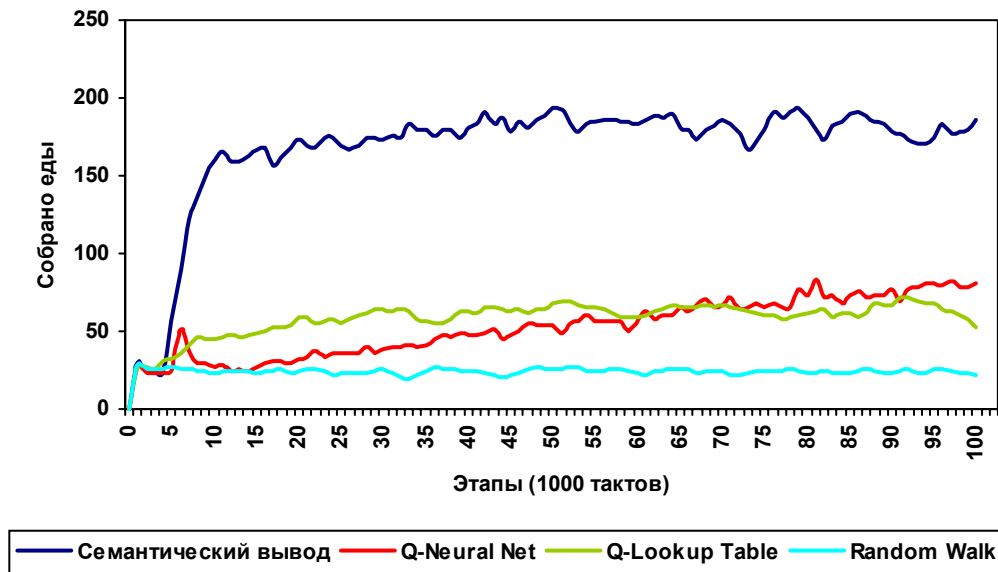


Рис. 6. Количество «еды», собранное аниматом с разными системами управления.

сплошной линией – подцели первого типа. Подцели, обозначенные двойной линией – это подцели второго типа, поскольку они увеличивают вероятность достижения вышестоящей цели, но при этом дальнейшие действия определяются неоднозначно. К примеру, если после достижения подцели «Еда на Севере» анимат направлен на восток, то ему надо повернуть налево, а если он направлен на запад – то направо.

Приведем пример, каким образом анимат использует выработанную им иерархию ФС для достижения цели. Допустим, что анимат ориентирован на восток, и сенсор «еда на северо-востоке» обнаружил еду. Очевидно, что в данной ситуации, чтобы получить результат, анимат должен совершить три действия: шагнуть вперед, повернуть налево и шагнуть вперед. Рассмотрим, какие цепочки запросов сформирует система управления аниматом, начиная с самой верхней системы:

*Еда здесь → Еда на Севере & Направление Север → Еда на Севере.*

Последняя функциональная система по закономерности

*(еда на северо-востоке)&(направление = восток)&(действие = шаг) → (еда на севере)*

с вероятностью 1 запустит на выполнение действие «шаг». Таким образом, будет достигнута подцель «Еда на Севере». По закономерности

*(еда на севере)&(направление = восток) &(действие = налево) →(еда на севере & направление = север)*

с вероятностью 1 анимат повернет налево и достигнет подцель «Еда на Севере & Направление Север». И, наконец, по закономерности

*(еда на севере) & (направление = север) & (действие = шаг) → (еда здесь)*

анимат шагнет вперед, в результате чего будет достигнута цель, и еда будет «съедена».

Для того чтобы оценить эффективность описанной модели, было решено провести сравнение данной системы управления с системами, построенными на основании теории обучения с подкреплением (Reinforcement Learning), описанной в работах [4].

Для сравнения мы выбрали две системы управления, построенные на основе популярного алгоритма обучения с подкреплением Q-Learning. Суть алгоритма заключается в последовательном уточнении оценок суммарной величины награды  $Q(s, A)$ , которую получит система, если в ситуации  $s_t$  она выполнит действие  $A_t$ , по формуле:

$$Q^{(i+1)}(s_t, A_t) = Q^{(i)}(s_t, A_t) + \alpha(r_t + \gamma \max_A Q^{(i)}(s_{t+1}, A) - Q^{(i)}(s_t, A_t)).$$

Первая из этих двух систем (Q-Lookup Table) основана на использовании таблицы, которая содержит Q-значения для всех возможных ситуаций и действий. Изначально эти значения таблицы заполняются случайным образом. В процессе работы в каждый такт времени система совершает действие и уточняет соответствующие Q-значения.

Вторая система (Q-Neural Net) использует аппроксимацию функции  $Q(s, A)$  при помощи нейронных сетей. При этом для каждого возможного действия  $A_i$  используется своя нейронная сеть  $NN_i$ . В каждый такт времени система выбирает действие, чья нейронная сеть выдаст наибольшую оценку Q-значения, после чего действие совершается и происходит адаптация весов соответствующей нейронной сети по алгоритму Back Propagation.

Для эксперимента было выбрано поле размером 25 на 25 клеток. Весь период функционирования анимата был разбит на этапы по 1000 шагов (тактов). Изначально все системы делают 5000 случайных шагов, чтобы накопить статистику и исследовать окружающую среду, после чего начинают функционировать в обычном режиме. Оценивалось, какое количество еды соберет анимат с разными системами управления за каждый этап работы. Очевидно, что после того как система управления полностью обучится и достигнет своего оптимального поведения, анимат начнет собирать примерно одно и то же количество еды за один этап. Таким образом, можно оценить как эффективность каждой системы управления в целом, так и скорость её обучения.

Было проведено несколько экспериментов, в ходе которых исследовались эффективность и скорость обучения каждой системы управления в случаях с различной плотностью заполнения среды едой.

Система управления на основе семантического вероятностного вывода во всех случаях показывала более высокую скорость обучения и качество работы по сравнению с системами Reinforcement Learning. Обычно уже после 8,000-10,000 тактов работы система полностью обучалась и достигала своего оптимального поведения.

Система управления на основе нейронных сетей (Q-Neural Net) с использованием отдельной сети для каждого действия не смогла обучиться даже после 100,000 шагов. По

этой причине в систему были внесены следующие изменения. Во-первых, было решено использовать по одной нейронной сети для каждого действия и каждого направления анимата. Таким образом, с учетом трех действий и четырех возможных направлений, система использовала 12 нейронных сетей. Во-вторых, задача для системы управления была упрощена путем уменьшения вдвое числа сенсоров, т.е. вместо 18-ти сенсоров на вход системы подавалось только 9 сенсоров, информирующих о наличии еды. Полученная таким образом система уже смогла научиться решать поставленную задачу.

Нейросетевая система управления оказалась чувствительна к плотности заполнения поля едой. При большой плотности заполнения система обучается быстрее. Так, при количестве еды на поле равным 200, система достигает оптимального поведения в среднем через 30,000-40,000 эпох, тогда как при количестве еды равно 100 – через 70,000 – 80,000. При количестве еды равным 50 и меньше система перестает обучаться. Вероятно, это связано с чувствительностью алгоритма Back Propagation к числу показов отдельных примеров - чем меньше еды на поле, тем меньшее количество раз нейронной сети предъявляются примеры, связанные с достижением аниматом цели, соответственно, тем меньше влияние этих примеров на процесс обучения. Кроме того, при уменьшении количества еды на поле, сильно возрастает нестабильность работы системы. Т.е. при небольшой плотности заполнения едой, анимат чаще попадает в ситуации, когда он длительное время не может встретить еду, в результате чего система начинает «разобучаться».

Системы управления с использованием таблицы Q-значений показала плохое качество работы, что объясняется достаточно большим количеством возможных ситуаций: с учетом трех действий и четырех направлений возможно 9984 различных ситуаций. Результаты экспериментов показывают, что даже после 100,000 тактов работы система управления в среднем просматривает только около 4000 ситуаций. Таким образом, даже после длительного времени обучения возникают ситуации, когда система реагирует неадекватно. Кроме того, система оказалась крайне нестабильной.

Система на основе таблицы значений, также как и система на основе нейронных сетей, оказалась чувствительной по отношению к плотности заполнения среды едой. Однако, в отличие от нейросетевой системы управления, данная система лучше обучается при небольшой плотности заполнения. Это связано с тем, что при малом количестве еды вероятность попасть в ситуацию, в которой анимат уже был ранее, значительно выше, чем при большом количестве еды. Соответственно, система управления обучается поведению в наиболее вероятных ситуациях. Однако, с уменьшением количества еды также возрастает и нестабильность системы. Наилучшие результаты были достигнуты при количестве еды равном 100.

На рис.6 представлены результаты экспериментов для случая, когда количество еды на поле поддерживалось равным 100. Для каждой системы управления рассчитывались средние значения по результатам 20-ти испытаний.

## ВЫВОДЫ

В данной модели анимата используется только простейший способ формирования подцелей по сравнению с теми, которые описаны в [8]. Но уже эта возможность дает значительные преимущества в обучении. Поэтому необходимо провести дальнейшие исследования возможности автоматического формирования целей, которые, в конечном итоге, должны дать возможность автоматического построения иерархии целей как это описано в Теории Движений Н.А.Бернштейна.

### Литература

1. Витяев Е.Е. Целеполагание как принцип работы мозга // Модели когнитивных процессов (Выч. сист., 158), Труды ИМ СО РАН, Новосибирск, 1997, с.9-52.
2. Витяев Е.Е. Принцип работы мозга и процесс познания в науке и искусстве. Изд. НГУ, Новосибирск, 1995, с.64
3. Витяев Е.Е. Вероятностное прогнозирование и предсказание как принцип работы мозга // Измерение и Модели Когнитивных Процессов (Выч. системы, 162), Новосибирск, 1998, Стр. 14-40.
4. Витяев Е.Е. Формальная модель работы мозга, основанная на принципе предсказания // Модели Когнитивных Процессов. (Выч. Сист., 164), Новосибирск, 1998, Стр. 3-61
5. Витяев Е.Е. Формализация двух принципов работы мозга // Нейроинформатика и ее приложения (Материалы VIII Всероссийского семинара 6-8 октября 2000, Красноярск), Под ред. А.Н.Горбаня, Красноярск, 2000, с.33-34.
6. Eugenie E. Vityaev, Irina V. Starikova. Two new Principles of Brain Activity. In: PAST CONFERENCES PROCEEDINGS SCI 2002/ISAS 2002 (The 6th World Multiconference on Systematics, Cybernetics and Information, July 14-18, 2002, Orlando, Florida) VOLUME XIII.
7. Михиенко Е.В., Витяев Е.Е. Моделирование работы функциональной системы, VI Всероссийская научно-техническая конференция «Нейроинформатика-2004». Сборник научных трудов. В 2-х частях. Ч.2., М.: МИФИ, 2004, 124-129.
8. Витяев Е.Е. Объяснение Теории Движений Н.А.Бернштейна. VII Всероссийская научно-техническая конференция «Нейроинформатика-2005», Сборник научных трудов, часть 1., Москва, 2005, стр.234-240



9. Витяев Е.Е. Логика работы мозга. Проблемы нейрокибернетики. т.2, (материалы XIV-ой Международной конференции по нейрокибернетике), Ростов-на-Дону, 2005, с.14-17.
10. Витяев Е.Е. Семантический подход к созданию баз знаний. Семантический вероятностный вывод. // Выч. системы, вып.146, Новосибирск, 1992, с. 19-49.
11. Sutton R. and Barto A. Reinforcement Learning: An Introduction. – Cambridge: MIT Press, 1998. See also: <http://www-anw.cs.umass.edu/~rich/book/the-book.html>
12. Витяев Е.Е. Метод обнаружения закономерностей и метод предсказания // Эмпирическое предсказание и распознавание образов (Выч. сист., вып. 67), Новосибирск, 1976, с. 54-68
13. Кендал М., Стьюарт А. Статистические выводы и связи // М., Наука, 1973, с.899.