



# B.A.R.: REPRESENTACIONES DISTRIBUIDAS PARA EL LÉXICO BILINGÜE

**OLGA SOLER VILAGELIU**  
Universitat Autònoma de Barcelona

## Resumen

La investigación sobre el acceso al léxico bilingüe no ha llegado a conclusiones definitivas a propósito de la organización de las entradas léxicas. Algunos efectos encontrados en trabajos recientes, como por ejemplo el efecto cognaticio, sugieren un almacenaje común para las entradas léxicas de los dos idiomas del hablante bilingüe, contradiciendo así aparentemente algunas hipótesis que postulan léxicos independientes y específicos para cada una de las lenguas. El Modelo Bilingüe para Representaciones de Acceso (*Bilingual Access Representations model*, BAR), basado en el modelo de Seidenberg y McClelland (*Model for Word Recognition and Naming*, 1989), es un modelo conexionista para el lexicón bilingüe que intenta simular esta hipótesis de almacenaje común, usando para ello representaciones distribuidas de las entradas léxicas. En este artículo se describe el modelo, su implementación y los resultados obtenidos en las simulaciones llevadas a cabo. Los análisis llevados a cabo después de las simulaciones con la red muestran la idoneidad de las representaciones distribuidas para modelar tanto el proceso del aprendizaje de nuevas entradas léxicas como los efectos empíricos de su similitud formal.

**Palabras clave:** Redes neurales, Acceso al léxico, Bilingüe

## Abstract

Research on bilingual lexical access has not yet arrived at definite conclusions on the organisation of lexical entries. Recent findings of cross-language effects, such as the cognate effect, seem to contradict previous hypotheses of language specific lexica, suggesting that some lexical entries of both languages might be stored together. The Bilingual Access Representations model (BAR), based on Seidenberg & McClelland's (1989) Model for Word Recognition and Naming, aims at capturing this feature of common storage by means of using distributed representations for lexical entries. The characteristics of the model, its implementation and the results obtained are presented in this paper. The analyses performed after the training show the suitability of distributed representations to model the process of learning lexical entries and the effects of their formal similarity.

**Key words:** Neural Networks, Lexical Access, Bilinguals.

\*El trabajo presentado en este artículo fue desarrollado en el Instituut voor Perceptie Onderzoek, (I.P.O.) en Eindhoven (Países Bajos). La autora quiere dar las gracias a Rudy van Hoe por su dirección, su colaboración y su apoyo; a Robert Hofsink por la implementación de la red y muchas buenas ideas; y a Don Bouwhuis por su confianza y especialmente por su ayuda en la elaboración de las primeras versiones del manuscrito. Hartelijk bedankt. Mi agradecimiento sincero también para Eduardo Navarrete y un revisor anónimo, quienes contribuyeron a dar los toques finales a esta versión.

Correspondencia: Olga Soler Vilageliu. Àrea Bàsica. Dept. Psicologia de l'Educació. Facultat de Psicologia- Edifici B. Universitat Autònoma de Barcelona. 08193- Bellaterra- Barcelona. correo electrònic: olga.soler@uab.es.

## Introducción

En este artículo presentaré un modelo para las representaciones de acceso al léxico bilingüe elaborado desde el marco teórico conexionista. Este modelo, basado en el modelo de Seidenberg y McClelland (1989), tiene como característica más importante que usa representaciones distribuidas para representar las entradas léxicas. La posibilidad de usar este tipo de representaciones se adoptó como una alternativa a los modelos clásicos de organización léxica bilingüe, los cuales no parecen explicar satisfactoriamente algunos de los efectos encontrados en la investigación. El modelo que se presenta no pretende resolver todos los problemas de representación del léxico bilingüe, sino explorar las posibilidades de un tipo de representación distinta para simular tanto la organización léxica bilingüe como el aprendizaje del vocabulario en una segunda lengua.

El trabajo revisa en primer lugar los distintos modelos sobre la organización del léxico bilingüe aparecidos en las últimas décadas, junto con el estudio de un factor que parece tener especial importancia en esta organización: la semejanza ortográfica. Seguidamente, se presenta el modelo Bilingüe para las Representaciones de Acceso (BAR) como una alternativa a las propuestas anteriores que permite integrar los efectos de esta variable. En el apartado siguiente (Método) se detallarán aspectos tanto de la implementación como de la simulación del modelo, con la descripción de la estructura de la red y del código empleado para describir la información fonética y ortográfica en la red. En el apartado Resultados se exponen los análisis llevados a cabo sobre los datos recogidos durante y después de las simulaciones. Estos análisis se centran en el aprendizaje cuantitativo de la red (medido por la tasa de error de la red); y en aprendizaje cualitativo (análisis de los tipos de errores realizados por la red). Finalmente, un análisis de conglomerados (*cluster analyses*), efectuado sobre los patrones de activación de las unidades internas de la red, permite estudiar cómo han sido representadas internamente las palabras. En las Conclusiones se señalan las ventajas e inconvenientes del modelo y se sugieren nuevas vías de investigación.

## Organización del Léxico Bilingüe: aspectos generales

Desde la década de 1950, cuando se elaboraron las primeras clasificaciones de hablantes bilingües, hasta la actualidad, se han propuesto distintos modelos para describir la organización del lexicón bilingüe<sup>1</sup>. La arquitectura de estos modelos es similar a la de los modelos clásicos del lexicón monolingüe (por ejemplo, Collins y Loftus (1975)), y por lo tanto constan de dos niveles de representación: el nivel semántico o conceptual y el nivel léxico. En opinión de De Groot (1992), estos modelos son excesivamente simples: para describirlos se usan diagramas de flujo<sup>2</sup> que no especifican de manera muy clara cómo se representa la información en los distintos niveles ni tampoco qué tipo de información se almacena en ellos.

Es necesario tener en cuenta que, hasta los años noventa, la investigación sobre la organización léxica y el acceso al léxico bilingüe se ocupó relativamente poco de cuestiones como la interacción entre la información fonológica y la ortográfica durante el reconocimiento de palabras. En consecuencia, no se incorporaron a los modelos propuestos versiones bilingües de los modelos de doble ruta o ruta única. Tampoco se especificó en esos modelos si las representaciones de las entradas léxicas en el nivel léxico son ortográficas o fonológicas o ambas cosas a la vez.

1 El lector interesado en una revisión de los distintos modelos puede consultar Soler Vilageliu, 1995.

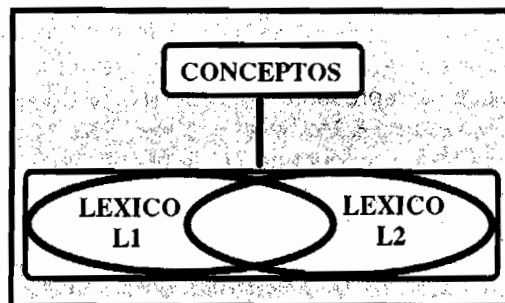
2 *Arrows and boxes* en el original.

Por otra parte, tampoco se precisa el tipo de información representada en el nivel conceptual o semántico. Como De Groot y Barry (1992) señalan, en la literatura sobre el lexicón bilingüe los términos "significado", "representaciones de significado" y "representaciones conceptuales" son intercambiables.

Los dos modelos de más influencia en los años ochenta constituyen un ejemplo de las características que se acaban de describir. Ambos modelos constan de una red de nodos de información para el nivel semántico y dos redes de nodos de información para el nivel léxico, es decir, una red léxica específica para cada lengua. El debate se centró en las relaciones entre estos dos léxicos específicos de lengua: la hipótesis de la Mediación por el Concepto (Potter, So, von Eckart & Feldman, 1984) supone que estas dos redes léxicas sólo se relacionan a través de sus respectivos vínculos con el nivel conceptual o semántico, mientras que la hipótesis de la Asociación de Palabras (Kirsner, Smith, Lockart, King & Jain, 1984) apoyaba asociaciones directas entre las representaciones léxicas dentro del mismo nivel. El foco de atención de estas investigaciones estaba más situado en describir la organización general del lexicón bilingüe que en estudiar aspectos del reconocimiento bilingüe de palabras como proceso cognitivo.

Por otra parte, la existencia de un léxico específico para cada una de las lenguas del hablante bilingüe ha sido cuestionada por varios investigadores recientemente (Alpitsis, 1990; Beauvillain, 1992; De Groot & Nas, 1991; García Albea, Bradley, Sánchez Casas & Forster, 1985; Sanchez-Casas, Davis & García-Albea, 1992). Estos autores proponen que el léxico debe ser parcialmente común para las dos lenguas y parcialmente específico para cada lengua, aunque defienden esta hipótesis con distintos argumentos.

Para Cécile Beauvillain (1992), el factor que organiza el léxico bilingüe no es la lengua a la cual pertenecen las palabras, sino la frecuencia con la que determinadas cadenas ortográficas aparecen en cada una de las lenguas. Después de realizar varios experimentos de reconocimiento visual de palabras con sujetos bilingües, Beauvillain sugirió que la representación interna de las palabras compuestas por segmentos no específicos de lengua (combinaciones de letras que pueden darse con igual frecuencia en las dos lenguas que el sujeto bilingüe conoce) se almacenan en la parte común del léxico bilingüe, mientras que las palabras con ortografía específica de cada lengua se almacenan en el léxico correspondiente a esta lengua<sup>3</sup>. Una representación gráfica del modelo de Beauvillain puede verse en la Figura 1.



**Figura 1.- Modelo de Cécile Beauvillain (1992). El nivel léxico está compuesto por dos lexicones específicos de idiomas superpuestos**

<sup>3</sup> Los ejemplos que Beauvillain (1992) usa son los siguientes: "CREAM", cuyas frecuencias de dígrafos suman 138 en Inglés y 52 en Francés, es una palabra específica del Inglés. En "TRADE", en cambio, las frecuencias de dígrafos suman 90 en Inglés y 96 en Francés, por lo cual no es una palabra específica para ninguna de las dos lenguas.

## Similitud Ortográfica

Pero si Beauvillain consideraba la frecuencia de los segmentos ortográficos en cada lengua el factor clave para la organización del léxico bilingüe, otros investigadores consideran que este factor debe ser una combinación de las variables ortografía y significado, aportando también evidencia empírica (Alpitsis, 1990; De Groot & Nas, 1991; García Albea, et al., 1985; Sanchez-Casas, et al., 1992). Estos autores han realizado experimentos utilizando como estímulos traducciones cognaticias, es decir, palabras en dos idiomas que son similares tanto en significado como en forma grafémica. Bajo el paradigma de la facilitación enmascarada (Forster, 1987), los resultados obtenidos en distintas tareas parecen indicar que sólo las palabras cognaticias facilitan el procesamiento de su traducción en la otra lengua, mientras que las traducciones sin ningún parecido ortográfico no se facilitan entre ellas. Este efecto se conoce como el efecto cognaticio (*cognate effect*). Tales resultados llevaron a los investigadores a concluir que las traducciones cognaticias (relacionadas tanto por su significado como por su forma ortográfica) están representadas en el nivel léxico por un mismo nodo de información. Las representaciones de estas entradas léxicas constituirían pues el área compartida por las dos lenguas en el léxico bilingüe, siendo la combinación de la similitud ortográfica y la similitud de significado la que produce el efecto cognaticio<sup>4</sup>.

## Génesis del Modelo Bilingüe para las Representaciones de Acceso

Revisadas estas recientes aportaciones al estudio de la organización léxica bilingüe, dos aspectos fueron los que sugirieron la conveniencia de ensayar un modelo conexionista usando representaciones distribuidas como el que presentamos aquí: en primer lugar, la propuesta de un léxico común para las dos lenguas; y en segundo lugar, la importancia de la ortografía como factor organizador del léxico bilingüe.

El desarrollo de un modelo de lexicón bilingüe que refleje estas características desde el enfoque clásico presenta algunos problemas (ver Soler Vilageliu, (1995) para una discusión más detallada). Como ya se ha visto, los modelos clásicos constan de redes compuestas por nodos de información, y cada uno de estos nodos corresponde a una entrada léxica. Establecer distintos tipos de relaciones entre estos nodos en función de su semejanza ortográfica o semántica o de su pertenencia a una u otra lengua es complicado.

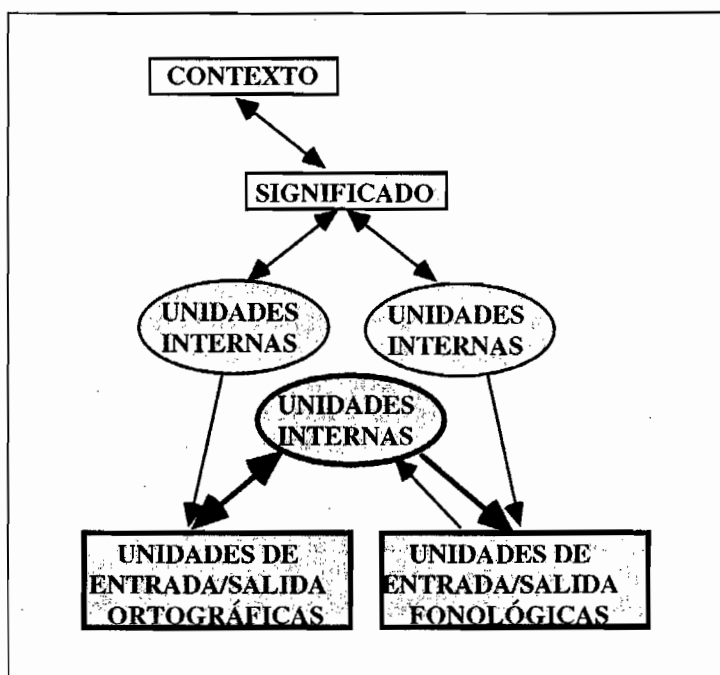
En cambio, la posibilidad que ofrecen los modelos conexionistas de representar las entradas léxicas de manera distribuida nos pareció una alternativa interesante. Usando una representación distribuida, cada entrada léxica se representa a través del patrón de activación de un número determinado de unidades (Seidenberg & McClelland, 1989). Este tipo de representación supone un cambio en el concepto de léxico, el cual no se compone de una red de representaciones de entradas léxicas sino de nodos con información subléxica. En consecuencia, la representación distribuida implica también un cambio en la forma de conceptualizar las relaciones entre las entradas léxicas: mientras que los modelos clásicos representan las relaciones de semejanza entre dos entradas léxicas a través de la proximidad de dos nodos en la red, un modelo con representaciones distribuidas representa la semejanza de dos entradas léxicas a través de la semejanza de los patrones de activación que las representan en el nivel léxico.

---

4 Podría argumentarse que, puesto que los experimentos citados utilizaron estímulos experimentales que estaban relacionados semánticamente (traducciones cognaticias o no cognaticias), el efecto específico de la ortografía similar nunca fue evaluado.

Debido a estas características, la representación distribuida de las entradas léxicas permite asumir de una manera estructural los dos aspectos mencionados: por un lado, las entradas léxicas de ambos idiomas pueden representarse en una misma entidad léxica; por el otro, las palabras compuestas por segmentos ortográficos similares se representarán por un patrón de activación similar.

Finalmente, un fenómeno no abordado por los modelos clásicos es el aprendizaje de nuevas palabras, aunque se supone que para incorporar nuevas entradas al léxico deben crearse nuevos nodos de representación. Un modelo de representación distribuida, en cambio, no necesitaría la creación de nuevos nodos, puesto que un número determinado de nodos pueden representar un gran número de entradas léxicas distintas mediante variaciones en el patrón de activación. En cualquier caso, estas hipótesis y la idoneidad de un modelo con representaciones distribuidas para el léxico bilingüe debían evaluarse mediante la implementación de una red neural. Tomando como referencia el Modelo para Reconocer y Nombrar Palabras (Seidenberg & McClelland, 1989) (Fig. 2) diseñamos un modelo para el léxico bilingüe.



**Figura 2.- Modelo para reconocer y nombrar palabras de Seidenberg y McClelland (1989). La parte que los autores implementaron es la representada en negrita (adaptado de Seidenberg y McClelland, 1989)**

Para simular el léxico bilingüe decidimos implementar solamente la parte del modelo que también implementaron Seidenberg y McClelland, es decir, el nivel léxico, sin incluir información semántica o conceptual.

El modelo, llamado Modelo Bilingüe para Representaciones de Acceso (*Bilingual Access Representations Model*, en adelante BAR), es un marco teórico para la organización de la memoria léxica bilingüe que fue concebido para describir el aprendizaje de nuevas palabras

durante la adquisición de una segunda lengua. Puesto que no tiene representaciones semánticas, BAR no es un modelo de traducción: BAR aprende diferentes conjuntos de correspondencias entre las representaciones ortográficas y fonológicas de dos lenguas distintas, el Neerlandés y el Inglés.

Las hipótesis que pretendíamos comprobar son las siguientes:

1. El modelo de representaciones distribuidas es un modelo capaz de representar el aprendizaje de palabras en dos lenguas (limitando el concepto de lengua a una correspondencia particular entre grafemas y fonemas) sin necesidad de almacenar las dos lenguas separadamente.

2. El aprendizaje de palabras en las dos lenguas se producirá de una manera plausible, es decir:

- A) Se producirá con sensibilidad a variables que han mostrado sus efectos repetidamente en estudios empíricos: la frecuencia de uso de las palabras en cada lengua y la longitud.
- B) Las palabras de similar estructura ortográfica y fonológica serán representadas también de forma similar.

La descripción de la red neural, el entrenamiento, los resultados obtenidos y las conclusiones a las que se llegaron para la última versión de BAR se exponen en los próximos apartados.

## Método

### BAR: Implementación (El Sujeto)

#### Arquitectura y Código

El modelo BAR tiene básicamente la misma arquitectura que la parte implementada del modelo de referencia (Seidenberg & McClelland, 1989), como puede verse en la Figura 3. La red consiste en dos conjuntos de unidades de entrada y salida, uno para la información ortográfica y uno para la información fonológica; más un conjunto de unidades internas que constituyen el nivel de representación léxica del modelo. Además, se ha introducido un nuevo conjunto de unidades de entrada, llamadas Unidades de Tarea. Puesto que, a diferencia del modelo de Seidenberg y McClelland, la información que BAR recibe es en dos lenguas, este grupo de unidades tienen la función de evitar lo que se conoce como Interferencia Catastrófica (McCloskey & Cohen, 1989). La interferencia catastrófica es el efecto por el cual la red puede olvidar completamente la información ya aprendida durante el proceso de adquirir nueva información (otras medidas fueron tomadas con el mismo objetivo, que se explicarán en el apartado Entrenamiento).

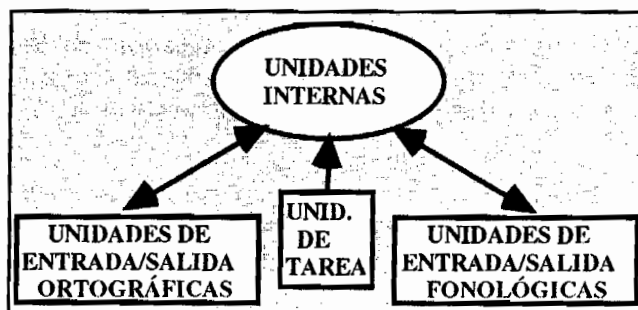


Figura 3.- Modelo Bilingüe para las Representaciones de Acceso (BAR)

Como en el modelo de Seidenberg y McClelland (1989), el algoritmo de aprendizaje usado en BAR es el algoritmo de retropropagación (*back propagation*). Durante el entrenamiento, este algoritmo permite ajustar los pesos de las conexiones entre unidades mediante la comparación de la activación de salida esperada (*target*) con la activación realmente obtenida como patrón de salida por la red (*output*).

El número de unidades que forman la red está condicionado por el tipo de código que se usa para proporcionar información a la red. Para BAR no se usó el código original empleado por Seidenberg y McClelland (1989), sino una adaptación del propuesto por McWhynney y Leinbach (1991), que fue usado también por Seidenberg en un modelo posterior (Daugherty & Seidenberg, 1994), con mejores resultados.

El código usado en BAR consiste en un grupo de unidades que codifican la posición y una característica mínima. Estas unidades se organizan en una plantilla silábica donde se representa cada palabra. La plantilla silábica tiene esta forma:

CCCVVCCCVVCCCVVCCCVVCCCVVCC

donde C significa Consonante y V Vocal. Esta plantilla puede codificar palabras de hasta 5 sílabas. Puesto que BAR necesita tanto información ortográfica como fonológica, el patrón de entrada contiene dos plantillas: una plantilla ortográfica y una plantilla fonológica. La plantilla ortográfica representa cada palabra según el alfabeto; y la plantilla fonológica representa la pronunciación de cada palabra mediante los fonemas que la componen.

El código ortográfico distingue entre 22 consonantes y 6 vocales. Puesto que el modelo se centraba en el desarrollo de patrones internos, no era necesario dotar a la red de las características visuales de los caracteres ortográficos. En su lugar, se usó un código binario para cada carácter. Según este código eran necesarias 3 unidades en cada posición para representar una vocal y 5 unidades en cada posición para representar una consonante.

El código fonológico distingue entre 32 consonantes y 39 vocales (incluyendo diptongos), es decir, un total de 71 fonemas. Una versión piloto de BAR mostró que, usando también un código binario arbitrario para cada fonema, la red sustituía unos fonemas por otros muy distintos. Aparentemente, este código generaba unas representaciones internas muy distintas de las "humanas". Para resolver este problema, BAR usó un código basado en características fonológicas como, por ejemplo, sonoridad y punto de articulación. Se esperaba que con esta modificación la red, en sus errores, substituiría fonemas por otros fonemas similares, en lugar de hacer cambios arbitrarios. Se usaron 9 unidades por posición para describir tanto las vocales como las consonantes fonológicamente.

El número de unidades necesarias para codificar una palabra puede calcularse a partir de esta información. Cada plantilla consta de 18 posiciones para consonante, y cada consonante necesita 5 unidades para su codificación ortográfica y 9 para su codificación fonológica; por lo tanto la red necesita  $[18 \cdot (5+9)=]$  252 unidades para codificar las consonantes. Para las vocales, el número puede calcularse de la misma manera: hay 10 posiciones para vocales, que necesitan 3 unidades para la codificación ortográfica y 9 para la codificación fonológica, es decir  $[10 \cdot (3+9)=]$  120 unidades para la codificación de las vocales. El total de unidades de entrada<sup>5</sup> y de salida necesarias es pues de  $252+120=372$ .

Por lo que se refiere a las unidades internas, su número fue determinado en función de los resultados obtenidos en las primeras versiones de BAR, y fue establecido en 110 (ver Hofsink, 1996, para detalles). El número total de unidades de la red es  $372+110+375=857$  unidades. La implementación de la red y los programas de codificación fueron elaborados por Robert Hofsink (Hofsink, 1996) en lenguaje de programación C++. La simulación fue llevada a cabo en una Sun Workstation.

---

5 A las unidades de entrada es necesario añadir las tres unidades de Tarea, cuyo fundamento se explica en el apartado Entrenamiento.

## Materiales

Los conjuntos de palabras usados para el entrenamiento y la fase de test o evaluación se obtuvieron de la base de datos CELEX (Burnage, 1990), que proporcionaba la información ortográfica, fonética y de frecuencia de uso necesaria para el entrenamiento de BAR.

El conjunto de entrenamiento neerlandés contenía 8074 palabras, y el inglés 1906. Esta desproporción en el tamaño de los dos conjuntos de palabras responde a la voluntad de que en la primera fase del entrenamiento la red aprendiera un léxico de tamaño considerable, como corresponde al de un hablante adulto, mientras que en la segunda fase de entrenamiento el número de palabras fuera equivalente al aprendizaje de la lengua inglesa por parte de un hablante neerlandés adulto<sup>6</sup>.

Debido a la limitación del número de unidades de entrada y salida de la red, la selección de palabras para los conjuntos de entrenamiento se hizo con el criterio de longitud: todas las palabras tenían entre 3 y 14 caracteres ortográficos. El segundo criterio de selección fue la frecuencia de las palabras, que se utilizaba, como en Seidenberg y McClelland (1989) para calcular el número de presentaciones de la palabra en cada época de entrenamiento (ver el próximo apartado Entrenamiento para la justificación teórica de esta característica). Para el conjunto de entrenamiento en neerlandés se seleccionaron palabras con una frecuencia mínima de 6 apariciones por millón y un máximo de 2100. Las palabras inglesas se escogieron de frecuencia media-alta, siguiendo el criterio según el cual las primeras palabras que se aprenden de un idioma son las más comunes y, por lo tanto, más frecuentes.

El formato de cada línea del archivo de datos obtenido de CELEX era el siguiente:

<ortografía de la palabra> <fonología de la palabra> <frecuencia por millón>

Ejemplo: example, lgz#mpP, 278

Un programa, diseñado por Robert Hofsink para esta función, transformaba los archivos directamente obtenidos de CELEX en los archivos de patrones que la red recibía para el entrenamiento.

La frecuencia no formaba parte del patrón de entrada administrado a la red, sino que se usaba para calcular la probabilidad de presentación de un determinado patrón en cada época de entrenamiento. Esta probabilidad se obtiene aplicando una transformación logarítmica al valor de la frecuencia por millón consignada en CELEX, usando la fórmula siguiente:

$$p = K \cdot \log(\text{frecuencia} + 2)$$

Seidenberg y McClelland (1989) sugirieron que el valor de la constante  $K$  debe ser escogido con referencia a la palabra más frecuente, cuya probabilidad de ser presentada a la red debe ser de  $p=0,93$ . Según este criterio, el valor de  $K$  fue establecido en 0,28.

El formato del archivo de patrones de entrenamiento era el siguiente:

<p-ortografía de la palabra-activación u. de tarea> <patrón de entrada> <patrón target>

El valor  $p$  determina las veces que una palabra aparece en el subconjunto de palabras que se presentan a la red en cada época de entrenamiento. La *ortografía de la palabra* no es usada por la red, sino que sirve solamente de identificación del patrón, y la *activación de las unidades de tarea* indica cuál debe ser esta activación, como se explica en el próximo apartado Entrenamiento. El *patrón de entrada* se compone de un patrón binario para la representación ortográfica y un patrón binario para la representación fonológica, como han sido descritos en el apartado anterior Arquitectura y Código. Finalmente, el *patrón target* u *objetivo* es igual al patrón de entrada, con la representación de los patrones binarios para las representaciones ortográficas y fonológicas. Durante el entrenamiento, la red comparará el patrón de activación de sus unidades de salida con el *patrón target* u objetivo, y corregirá sus respuestas aplicando la fórmula de aprendizaje de retropropagación.

<sup>6</sup> Evidentemente, con las limitaciones impuestas por el modelo, que sólo puede representar la información ortográfica y fonológica de las palabras.



Finalmente, para evaluar el aprendizaje de la red se usaron también como conjuntos de *test* un conjunto de no-palabras y un conjunto de palabras no entrenadas. El primer conjunto estaba formado por patrones de no-palabras, que se obtuvieron invirtiendo 1000 palabras del neerlandés. Las no-palabras resultantes no tenían segmentos orto-fonológicos naturales en neerlandés. El segundo conjunto estaba simplemente formado por 1000 palabras neerlandesas que no habían sido entrenadas y se presentaban por primera vez a la red.

### Entrenamiento (Procedimiento)

El entrenamiento de BAR se realizó en dos fases: en la primera fase, Entrenamiento Monolingüe, la red era entrenada con el conjunto de palabras en Neerlandés; y en la segunda, Entrenamiento Bilingüe, la red era entrenada con palabras en Inglés, además del conjunto original de palabras en Neerlandés. En una tercera fase de evaluación o fase de test, se presentaban a la red los conjuntos de no-palabras y de palabras no entrenadas, descritas en el anterior apartado. Es necesario señalar que en la fase de test las palabras se presentan una sola vez a la red con el fin de obtener la activación de sus unidades de salida. Esta activación se compara con el patrón de entrada, pudiendo realizarse de esta manera la evaluación del aprendizaje realizado.

Hay varios aspectos importantes en el entrenamiento de la red que se detallan a continuación.

En primer lugar, y como ya se ha mencionado, uno de los problemas que las redes neurales con representaciones distribuidas pueden presentar en la simulación de procesos de aprendizaje es la Interferencia Catastrófica (McCloskey & Cohen, 1989). Los aprendizajes realizados de forma serial implican que la información aprendida en primer lugar desaparezca por completo o parcialmente al adquirir nueva información. En el modelo que nos ocupa, el aprendizaje de palabras de la segunda lengua ocasionaría el olvido de las palabras de la primera lengua, efecto que no es deseable porque en la vida real no se produce.

Hay varias propuestas que sirven para reducir este efecto. McCloskey y Cohen (1989) recogen la propuesta de Rumelhart de añadir unidades de contexto a las unidades de entrada. Las unidades de contexto informan a la red de cuándo el contexto ha cambiado y por lo tanto la información nueva debe almacenarse junto con la del contexto anterior. Siguiendo esta propuesta, se incorporaron a BAR las Unidades de Tarea<sup>7</sup>.

Otra sugerencia es la presentada por Murre (Murre, 1993), quien atribuye al método de presentación de los patrones de entrenamiento la causa principal de la interferencia catastrófica. Su propuesta es utilizar un método aleatorio de entrenamiento que presente los patrones antiguos juntamente con los nuevos. De esta manera, los patrones nuevos son aprendidos y los viejos se siguen entrenando. Este método fue también utilizado en BAR durante la segunda fase de entrenamiento, cuando se presentaron aleatoriamente en cada época de entrenamiento patrones de palabras inglesas nuevos y patrones de palabras neerlandeses ya entrenados en la primera fase.

En segundo lugar, otro aspecto muy importante para el entrenamiento (que ya se ha introducido en el apartado anterior) es el tratamiento de la frecuencia de las palabras. Como es bien

---

<sup>7</sup> Las unidades de Tarea sirven a la red para identificar la lengua de las palabras que está recibiendo. La identificación de la lengua ha merecido especial atención por parte de algunos autores (Grainger & Dijkstra, 1992; Thomas & Plunkett, 1995), pero la discusión sobre este aspecto queda fuera de los límites de este artículo. Sin embargo, considero importante precisar que en BAR se incluyeron sólo 3 Unidades de Tarea para que no tuvieran una gran influencia en la elaboración de la representación interna de las palabras, puesto que, una vez alcanzado cierto nivel de aprendizaje, la estructura ortográfica y fonológica de las palabras debería ser suficiente para reconocer su pertenencia a una u otra lengua.

conocido, la frecuencia de las palabras es trascendental en el aprendizaje de la lengua, puesto que las palabras más frecuentes se aprenden con más rapidez. Por lo tanto, era importante reflejar esta variable en el entrenamiento de BAR. Este objetivo se cumplía, como se ha visto en Material, haciendo que la probabilidad de presentación de los patrones de palabra en cada época de entrenamiento fuera calculada en función de la frecuencia de la palabra.

El coeficiente de aprendizaje de la red fue fijado en 0.02 para evitar fluctuaciones amplias en la curva de aprendizaje. Este coeficiente relativamente bajo (en versiones anteriores de BAR se utilizó un coeficiente de 0.05), es más adecuado cuando se dispone de un conjunto de datos complejos muy amplio (Hofstink, 1996). Tanto la fase monolingüe como la fase bilingüe se llevaron a cabo durante 1000 épocas de entrenamiento, ya que la primera implementación de BAR mostró que este número era suficiente (Soler Vilageliu, 1996). En la tercera fase, la fase de evaluación, los nuevos patrones se presentaban a la red una sola vez para poder evaluar el aprendizaje alcanzado por la red. Los resultados obtenidos y el análisis de los mismos se presentan en el próximo apartado.

## BAR: Resultados

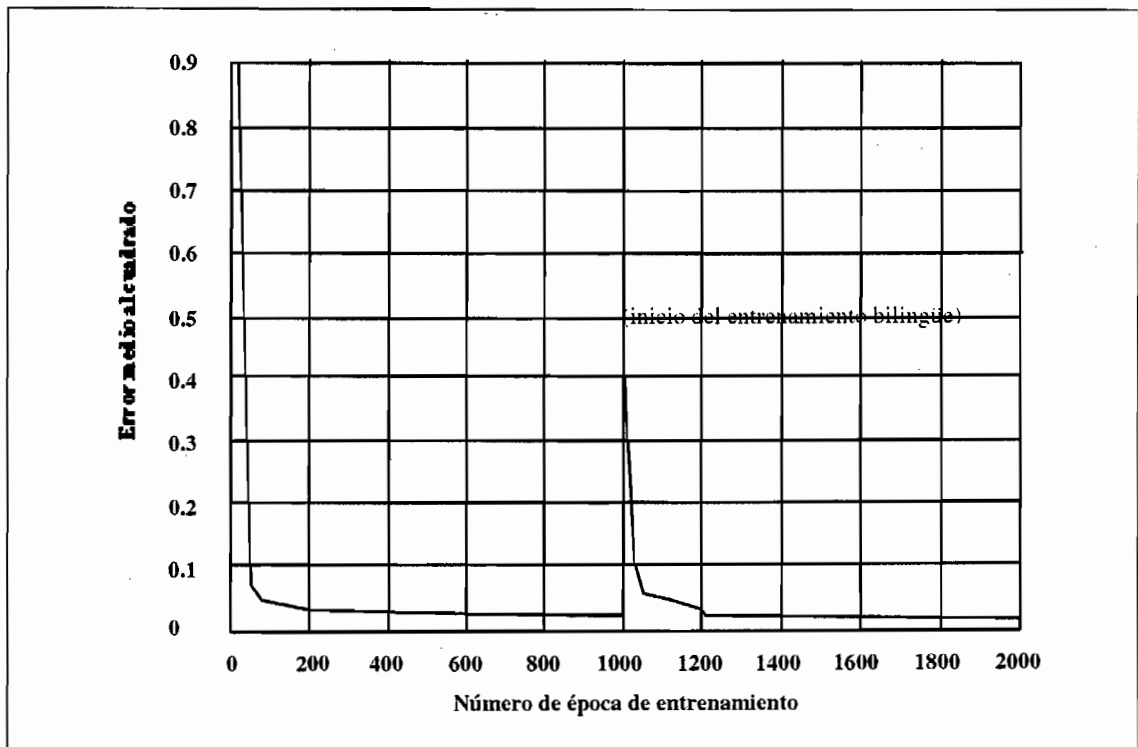
La evaluación de los resultados de BAR se llevó a cabo por distintos procedimientos. Por un lado, la tasa de error obtenida durante el entrenamiento (exactamente cada 200 épocas) permite la valoración cuantitativa del aprendizaje de la red; por el otro, la observación de los patrones de salida obtenidos después del entrenamiento permite la evaluación cualitativa, utilizando como variable el porcentaje de errores. Estos dos análisis permiten evaluar la eficacia del modelo para el aprendizaje de dos conjuntos de palabras de dos idiomas, y específicamente, de las relaciones ortografía-fonología de estos dos idiomas.

Pero además, una de las hipótesis planteadas respecto a la idoneidad de las representaciones distribuidas era que las palabras similares se representarían por patrones de activación similares en el nivel léxico, con independencia de su lengua. Para verificar esta hipótesis se aplicó un tercer análisis a las representaciones internas de BAR, que evaluó el parecido de los patrones de activación de las unidades internas mediante la técnica del análisis de conglomerados (*clustering*).

Finalmente, la red fue confrontada con dos conjuntos de palabras no entrenados: uno, de palabras neerlandesas que no pertenecían al conjunto de entrenamiento; y el segundo, con un conjunto de palabras neerlandesas invertidas, es decir, sin las características de lengua. En ambos casos se evaluó el porcentaje de errores.

## Curva de aprendizaje

En primer lugar, podemos observar el rendimiento de la red mediante la tasa de error (Figura 4). Esta tasa de error corresponde a la media al cuadrado del error de la red, que se calcula mediante la comparación entre la activación obtenida en las unidades de salida y la indicada en el patrón *target* u objetivo, según el algoritmo de aprendizaje retropropagación. La tasa de error puede registrarse durante el entrenamiento, sin necesidad de detenerlo, y se registró cada 200 épocas de entrenamiento. Tanto en el entrenamiento monolingüe como en el bilingüe esta tasa desciende bruscamente durante las primeras cien épocas de entrenamiento, para converger en un valor aproximado de 0.01. Puesto que el valor es similar después de ambas fases de entrenamiento, es posible concluir que el entrenamiento bilingüe no ha producido interferencia en los patrones aprendidos en la fase monolingüe.



**Figura 4.- Curva de aprendizaje de BAR durante las dos fases de entrenamiento**

### Aprendizaje cualitativo

El aprendizaje cualitativo de la red fue evaluado a partir de los errores en las representaciones de salida de la red, considerando tres tipos de errores posibles: en la representación ortográfica (sólo ortografía), en la representación fonológica (solo fonología) y en ambas representaciones (ortografía+fonología). La Tabla 1 muestra los porcentajes de palabras incorrectas obtenidos después de las dos fases de entrenamiento tanto para el neerlandés como para el inglés.

**Tabla 1.- Porcentajes de palabras incorrectas en cada idioma para las dos fases de entrenamiento**

Tipo de Error	Entrenamiento Monolingüe		Entrenamiento Bilingüe	
	palabras neerlandesas incorrectas (%)	palabras inglesas incorrectas (%)	palabras neerlandesas incorrectas (%)	Palabras inglesas incorrectas (%)
sólo ortografía	0.17	3.93	0.10	0.05
sólo fonología	0.36	40.03	0.17	1.00
ortografía+fonología	0.05	40.77	0.02	0.05
<b>Total</b>	<b>0.58</b>	<b>84.73</b>	<b>0.29</b>	<b>1.10</b>

Después del entrenamiento monolingüe, BAR tiene un rendimiento casi perfecto en las palabras neerlandesas (0.58% de error). Por otra parte, el modelo muestra un desconocimiento total de las palabras inglesas, que le fueron presentadas en fase de *test* (84,73% de errores)<sup>8</sup>. Evaluar el conocimiento de las palabras inglesas anterior al entrenamiento bilingüe es útil para contrastarlo con los resultados posteriores. Los resultados obtenidos después del entrenamiento bilingüe muestran que no ha habido ninguna interferencia de los patrones ingleses sobre los neerlandeses, puesto que el rendimiento en las palabras neerlandesas incluso ha mejorado, disminuyendo el porcentaje de errores (0,29%). El rendimiento en las palabras inglesas es también muy bueno, con sólo un 1,10% de errores sobre el total de palabras aprendidas.

### Frecuencia de aparición de las palabras

Dos variables fueron especialmente evaluadas después del entrenamiento: la frecuencia y la longitud de las palabras. Respecto a la frecuencia, es interesante verificar si las palabras menos frecuentes se han aprendido con igual corrección que las más frecuentes para los dos idiomas. Para estudiar el efecto de la frecuencia se seleccionó un subconjunto con las palabras de alta frecuencia (entre 200 y 1370 por millón, aproximadamente 250 palabras) y un subconjunto con las palabras de baja frecuencia (6 por millón, aproximadamente 250 palabras) tanto neerlandesas como inglesas. Para cada uno de estos subconjuntos se calcularon los porcentajes de palabras incorrectas igual como se hizo con el total de palabras. Estos porcentajes pueden observarse en la Tabla 2.

**Tabla 2.- Porcentajes de palabras incorrectas de alta y baja frecuencia en cada idioma para las dos fases de entrenamiento**

Tipo de Error	Entr. monolingüe		Entrenamiento bilingüe			
	Neerlandés		Neerlandés		Inglés	
	Alta frec.(%)	Baja frec.(%)	Alta frec.(%)	Baja frec.(%)	Alta frec.(%)	Baja frec.(%)
sólo ortografía	0	0.73	0	0	0	0
sólo fonología	0.28	2.20	0	1.10	0	6.00
ortografía+fonología	0	0	0	0	0	0
<b>Total</b>	<b>0.28</b>	<b>2.93</b>	<b>0</b>	<b>1.10</b>	<b>0</b>	<b>6.00</b>

Como se ve en la Tabla 2, las palabras de baja frecuencia tienen una desventaja mínima con respecto a las palabras de alta frecuencia. Este resultado responde a lo esperado: las diferencias en el aprendizaje según la frecuencia de las palabras deben reflejarse durante éste, pero deben ser mínimas al final del mismo. Sin embargo, los resultados para las palabras en

<sup>8</sup> Es interesante destacar que la mayoría de los errores en las palabras inglesas antes de la fase de entrenamiento bilingüe es debida a errores en la plantilla de representación fonológica (80,80%), mientras que los errores en la plantilla de representación ortográfica contribuyen en la mitad (44,70%) al porcentaje total. Esta diferencia es probablemente debida a que el código empleado para la representación fonológica se basa en las características articulatorias de los fonemas, que son distintas en los dos idiomas. Por el contrario, el código ortográfico tiene la misma descripción.

inglés son extremadamente buenos. Esto se debe a que el conjunto seleccionado de palabras inglesas tenía en promedio una frecuencia más alta que el de palabras neerlandesas, y por lo tanto la probabilidad de aparición de las palabras inglesas de baja frecuencia durante el entrenamiento bilingüe era más alta que la probabilidad de aparición de las palabras holandesas de baja frecuencia<sup>9</sup>. Es remarcable el hecho de que el porcentaje de palabras incorrectas más alto corresponde a la representación fonológica de las palabras inglesas de baja frecuencia, resultado que puede ser una combinación de dos variables: la interferencia de los patrones neerlandeses, y la propia irregularidad del Inglés en cuanto a su pronunciación.

### Longitud de palabra

La longitud de las palabras es la otra variable a observar, puesto que BAR usa un código con una plantilla de longitud fija, mientras que las palabras entrenadas son de longitud variable. Para evaluar el efecto de esta variable se seleccionaron dos subconjuntos de palabras para cada idioma, uno formado por las palabras más cortas y otro formado por las palabras más largas. Los cuatro subconjuntos estaban compuestos también por un número aproximado de 250 palabras. Los porcentajes de palabras incorrectas de cada subconjunto pueden verse en la Tabla 3.

**Tabla 3.- Porcentajes de palabras incorrectas largas y cortas en cada idioma para las dos fases de entrenamiento**

Tipo de Error	Entr. monolingüe		Entrenamiento bilingüe			
	Neerlandés		Neerlandés		Inglés	
	Palabras cortas (%)	palabras largas (%)	palabras cortas (%)	palabras largas (%)	palabras cortas (%)	palabras largas (%)
sólo ortografía	0	3.19	0	2.39	0	0.47
sólo fonología	0	4.38	0	2.39	0	6.98
ortografía+fonología	0	1.59	0	0.80	0	0.47
<b>Total</b>	<b>0</b>	<b>9.16</b>	<b>0</b>	<b>5.58</b>	<b>0</b>	<b>7.91</b>

Es evidente que la variable longitud de palabra es importante para el aprendizaje de la red: las palabras largas tienen el mayor porcentaje de error en ambas lenguas. La explicación para este efecto se encuentra en el sistema de codificación por plantillas. En estas plantillas las palabras se codifican de izquierda a derecha, de manera que las primeras posiciones están siempre ocupadas aunque la palabra sea muy corta, mientras que las posiciones de la derecha sólo se llenan cuando la palabra es larga. Puesto que las palabras largas sólo constituyen, aproximadamente, un 3% del total de palabras en cada conjunto de entrenamiento, una parte importante de conexiones de la red no se entrenan lo suficiente para poder aprenderlas.

<sup>9</sup> Este es un aspecto a considerar para nuevas versiones de BAR. Aunque las palabras que un estudiante de Inglés como segunda lengua aprende en primer lugar son seguramente las más frecuentes en este idioma, la ocurrencia real de estas palabras en el entorno lingüístico del estudiante es mucho menor que en el conjunto del idioma Inglés. En consecuencia, debería reducirse la probabilidad de presentación de las palabras de la segunda lengua.

## Representaciones Internas: Análisis de conglomerados

Como se ha introducido al inicio de este apartado, las representaciones internas de BAR se evaluaron mediante la técnica del análisis de conglomerados (*clustering*). Esta evaluación tenía como objetivo verificar que las palabras ortográficamente y fonológicamente similares se representan en BAR por patrones de activación similares. Aunque esta técnica supone una aproximación relativa, puesto que depende del número de palabras incluido en el conjunto, es útil para valorar la semejanza de los patrones de activación.

El conjunto seleccionado para el análisis estaba compuesto por 85 palabras<sup>10</sup>, incluyendo palabras similares dentro de cada idioma (*hair hear, broek broer*) y palabras similares entre los dos idiomas, tanto traducciones cognaticias (*hel hell*) como palabras que no tenían relación semántica (*vorm warm*), aunque esta distinción es irrelevante para BAR porque no tiene representación de significado. Los conglomerados (*clusters*) obtenidos confirmaron las expectativas, agrupando palabras por su parecido formal. Los siguientes ejemplos muestran los *clusters* obtenidos por debajo de la distancia media:

(((*hart hert*) *haar*) ((*hel hell*) *heel*)))  
 (*hair hear*)  
 (*boat coat*)  
 (*broek broer*)  
 (*bal ball*)

Aunque en estos *clusters* encontramos palabras de los dos idiomas también es importante señalar que los primeros que se forman son entre palabras de un mismo idioma (*hart hert, boat coat*). Esto parece indicar que las relaciones ortografía/fonología que se establecen en cada lengua pueden ser información suficiente para identificar su pertenencia a esta lengua.

## Evaluación con palabras no entrenadas y no-palabras

Para comprobar que BAR es capaz de capturar las características ortográficas y fonológicas de una lengua, la red se evaluó con dos conjuntos más de patrones (Hofsink, 1996). Como se ha descrito en el apartado Materiales, estos patrones fueron extraídos también de la base de datos CELEX. El cómputo de los errores cometidos por la red en estos dos conjuntos muestra su capacidad para representarlas correctamente, lo que puede tomarse como una medida de "reconocimiento". Los porcentajes obtenidos en estos dos conjuntos pueden observarse en la Tabla 4.

Tabla 4.- Porcentaje de palabras incorrectas para no-palabras y palabras nuevas

Tipo de Error	no-palabras incorrectas (%)	nuevas palabras incorrectas (%)
sólo ortografía	5.5	3.2
sólo fonología	22.8	5.1
ortografía+fonología	35.8	2.4
<b>Total</b>	<b>64.1</b>	<b>10.7</b>

<sup>10</sup> Como le parecerá obvio al lector, el análisis de conglomerados debe realizarse con un subconjunto de palabras, puesto que realizarlo con el conjunto total de palabras (unas 10000) sería imposible y, por ende, poco informativo.

La diferencia entre el porcentaje de error total para cada uno de los conjuntos muestra claramente que la red es capaz de representar correctamente la mayoría de palabras nuevas (10,7% de errores, 89,3% de palabras correctas), mientras que las no-palabras representadas correctamente son menos de la mitad (35,9%). Según Hofsink (1996) el rendimiento para las palabras nuevas sería mejor si la red tuviera menos unidades internas, puesto que esto forzaría a la red a usar más la redundancia ortográfica y fonológica del lenguaje natural. En cualquier caso, estos porcentajes indican la capacidad de la red para reconocer las restricciones fonotácticas del neerlandés.

## Conclusiones

Los resultados obtenidos en la simulación del modelo BAR muestran que éste es adecuado para representar el lexicón bilingüe. Como se decía en el inicio de este artículo, el interés de la implementación y evaluación del modelo residía básicamente en evaluar la idoneidad de la representación distribuida de las entradas léxicas en el léxico bilingüe. Este enfoque permite incorporar la noción de lexicón común y al mismo tiempo específico para las dos lenguas, no mediante las relaciones entre las entradas léxicas como en los modelos clásicos de organización léxica bilingüe, sino por la estructura propia de este modelo conexionista.

Como se ha mostrado en el apartado anterior, los resultados muestran en primer lugar que el modelo BAR es capaz de aprender un vocabulario extenso de palabras en dos idiomas (un total de 10000 palabras), aprendiendo y distinguiendo las distintas relaciones ortografía-fonología de cada una de ellas. Esto nos permite verificar la primera hipótesis expuesta en la Introducción.

Por otra parte, gracias al método de entrenamiento de la red, se ha simulado el efecto de la frecuencia de la palabra en el aprendizaje de los dos vocabularios. La red muestra también un efecto de la variable longitud de palabra, efecto que si bien no es del todo indeseable, no es un efecto esperado sino debido a limitaciones del código usado. El aprendizaje de las correspondencias ortografía-fonología en los dos idiomas es suficiente para su generalización a palabras nuevas que no se habían entrenado, como se ha podido comprobar para un conjunto de palabras neerlandesas. También se ha comprobado que la red representa incorrectamente las no-palabras que no tienen una estructura ortográfica y fonológica propia de estos dos idiomas. Estos aspectos confirman que la red lleva a cabo su aprendizaje de manera sensible a las variables establecidas de frecuencia y longitud, con cierta plausibilidad psicológica.

Estos resultados, junto con el análisis mediante la técnica del análisis de conglomerados realizado sobre los patrones de activación de las unidades internas, nos permiten suponer que el modelo BAR es capaz de simular los hallazgos de Beauvillain (1992) respecto a la representación de segmentos ortográficos no específicos de lengua, y también los del grupo de investigadores del efecto cognaticio (Alpitsis, 1990; De Groot & Nas, 1991; García Albea, et al., 1985; Sanchez-Casas, et al., 1992), respecto a la representación de las traducciones cognaticias, puesto que las traducciones cognaticias son similares formalmente.

Es evidente que BAR, al no contar con un nivel de representación semántico o conceptual, no es un modelo apto para describir la organización léxica bilingüe en su totalidad, pero ha mostrado su utilidad para representar el nivel de codificación de acceso al léxico y algunos efectos encontrados en la investigación empírica, como por ejemplo el efecto cognaticio descrito en la Introducción.

Nuevas investigaciones podrán probar la eficacia de las representaciones distribuidas para simular la organización léxica, tanto bilingüe como monolingüe, como han sugerido en recientes investigaciones sobre el efecto de vecindad otros autores (Ziegler & Perry, 1998). El modelo BAR constituye solamente una primera propuesta, desde un marco teórico diferente a los explorados hasta el momento, para la representación de la organización léxica bilingüe.

## Referencias

- Alpitsis, R. (1990). *Lexical Representation in Greek/English Bilinguals*. Honours, Monash University.
- Beauvillain, C. (1992). Orthographic and Lexical Constraints in Bilingual Word Recognition. In R.J. Harris (Ed.), *Cognitive Processing in Bilinguals* (pp. 221-236). Amsterdam: North Holland.
- Burnage, G. (1990). *Celex- A Guide for Users*. Nijmegen: Celex- Centre for Lexical Information.
- Collins, A.M. & Loftus, E.F. (1975) A Spreading-Activation Theory of Semantic Activation, *Psychological Review*, 82 (6), 407-428.
- Daugherty, K.G., & Seidenberg, M.S. (1994). Beyond rules and exceptions. In S.D. Lima, R.L. Corrigan, & G.K. Iverson (Eds.), *The Reality of Linguistics Rules* (pp. 353-388). Amsterdam/Philadelphia: John Benjamins.
- De Groot, A.M.B. (1992). Bilingual representations: A closer look at Conceptual Representations. In Frost & Katz (Eds.), *Orthography, Phonology, Morphology, and Meaning* Amsterdam: Elsevier.
- De Groot, A.M.B., & Barry, C. (1992) The Multilingual community: Introduction. *European Journal of Cognitive Psychology*, 4(4), 241-252.
- De Groot, A.M.B., & Nas, G.L.J (1991). Lexical Representation of Cognates and Noncognates in Compound Bilinguals. *Journal of Memory and Language*, 30, 90-123.
- Forster, K.I. (1987). Form-Priming with Masked Primes: The Best Match Hypothesis. In M. Coltheart (Ed.), *Attention and Performance XII. The Psychology of Reading*. London: Lawrence Erlbaum Ass.
- García Albea, J.E., Bradley, D.C., Sánchez Casas, R.M., & Forster, K.I. (1985). Cross Language priming effects in bilingual word recognition. In *Fifth Australian Language and Speech Conference*, Parkville, November:
- Hofsink, R.D.J.H. (1996). *Language modelling with artificial neural networks*. Eindhoven: IPO Rapport no. 1119.
- Kirsner, K., Smith, M.C., Lockart, R.S., King, M.L., & Jain, M. (1984). The Bilingual Lexicon: Language-Specific Units in an Integrated Network. *Journal of Verbal Learning and Verbal Behavior*, 23, 519-539.
- MacWhinney, B. & Leinbach, J. (1991) Implementations are not conceptualizations: Revising the Verb Learning Model. *Cognition*, 40, 121-157.
- McCloskey, M., & Cohen, N.J. (1989). Catastrophic interference in connectionist networks: The sequential learning problem. *The Psychology of Learning and Motivation*, 24, 109-165.
- Murre, J.M.J. (1993). The effects of pattern presentation on interference in back-propagation networks. In *Proceedings of the 14th Annual Conference of the Cognitive Science Society* (pp. 54-59).
- Potter, M.C., So, K.F., von Eckart, B., & Feldman, L.B. (1984). Lexical and Conceptual Representation in Beginning and Proficient Bilinguals. *Journal of Verbal Learning and Verbal Behavior*, 23, 23-38.
- Sanchez-Casas, R.M., Davis, C.W., & Garcia-Albea, J.E. (1992). Bilingual Lexical Processing: Exploring the Cognate/Non-Cognate Distinction. *European Journal of Experimental Psychology*, 4(4), 293-310.
- Seidenberg, M.S., & McClelland, J.L. (1989). A Distributed, Developmental Model of Word Recognition and Naming. *Psychological Review*, 96(4), 523-568.
- Soler Vilageliu, Olga (1995). Estudio Experimental del Bilingüismo: Revisión Histórica. *Anuario de Psicología*, 66, 19-36.
- Soler Vilageliu, Olga (1996). *Bilingual lexical access: A connectionist model*. Tesis Doctoral, Universitat Autònoma de Barcelona.
- Ziegler, J.C. & Perry, C. (1998) No more problems in Coltheart's neighborhood: resolving neighborhood conflicts in the lexical decision task. *Cognition*, 68, B53-B62.