

## Tanulmány

Sass Bálint

### Az igei szerkezetek algebrai struktúrája, avagy a duplakocka modell

#### Abstract

In this paper, firstly (1) we define the general „construction” concept, and describe an abstract model which captures this concept. This model is built on the notion of algebraic lattice. Secondly, (2) we apply the model to verbal expressions like *he takes their opinions into consideration* or *the team would get rid of its beloved mascot*. Then (3) by combining verbal expressions we make the model suitable for representing verbal expressions of a whole corpus together instead that of just one clause. Finally, (4) we outline the algorithmic framework, by which the model can become capable of identifying the so called real verbal expressions like *take something into consideration* or *get rid of something* in a corpus-driven way.

*Keywords:* lattice, algebraic model, verbal expression, lexical acquisition, corpus driven

#### 1 Bevezetés

*Hiszek abban, hogy ha adsz neki egy apró ajándékot, akkor újra kisüt a nap és pontot tehetsz az egész dolog végére. Az a gond, hogy részt veszel a vitában, de nem veszel részt a munkában. Legalább egy pillantást vess rá, hogy ne vethessék a szemedre a nemtörődömséget.*

Mondatainkban az igék mellett lehetnek „helyek”, egy hely (*hisz vmiben*) vagy több (*ad vkinek vmit*), és ezeket a helyeket a szövegben konkrét szavak „töltik ki” (*hisz abban és ad neki ajándékot*). Sok esetben a kitöltő szavak a konstrukció jelentésének megmaradása mellett viszonylag szabadon variálhatók, máskor azonban nem, mert a változtatás az eredeti jelentés elvesztésével jár. Egy ige mellett lehet egy (*kisüt a nap*) vagy több (*pontot tesz a végére*) ilyen utóbbi értelemben kötött hely. Olyan is van, hogy adott ige mellett vegyesen mindkét típus megjelenjen: a *részt vesz vmiben* esetén a tárgy fix, a *vmiben* helyen viszont sokféle szó válthatja egymást (*vitában, munkában* stb.). Sőt azzal is találkozunk, hogy egy ige melletti két hely közül egyszer az egyik van kötött módon kitöltve, másszor viszont a másik (*vet pillantást vmire és vet szemére vmit*).

A most bemutatott példák mind igei szerkezetek. Ebben a tanulmányban az igei szerkezeteknek egy új, formális, matematikai (algebrai) modelljét mutatjuk be. Ez a modell az igei szerkezetek egymáshoz való viszonyait is megragadja. Kezeli, magában foglalja a fent vázolt eseteket, beleértve az olyanfajta átfedéseket is, mint amiket a két *vet*-es szerkezet példáz. A modell alkalmasnak látszik arra, hogy haszonnal alkalmazzuk a különféle igei szerkezetek kutatásában, és hogy segítségével új, hatékony igeiszerkezet-kinyerő eljárásokat alkossunk.

Ahogy annak idején Kalmár (1964: 12) megfogalmazta, a matematikai nyelvészet „a matematikus tapasztalatait igyekszik gyümölcsöztetni a nyelv strukturális vizsgálata számára”. Jelen tanulmány így a klasszikus értelemben vett matematikai nyelvészet keretébe illik bele,

amennyiben matematikai struktúrákat alkalmaz nyelvi jelenségek leírása, modellezése, vizsgálata során. Az említett tanulmány többféle algebrai struktúrát említ, a hálókat nem. Most viszont éppen a hálók lesznek a középpontban, többféle hálóval is találkozni fogunk.

Az első fejezetekben részletesen bemutatjuk a duplakocka modell absztrakt formáját (2–4. rész, 13. oldal), aztán alkalmazzuk a modellt az igei szerkezetek reprezentálására (5. rész, 22. oldal), alkalmassá tesszük arra, hogy egy korpuszban fellelhető összes igei szerkezetet egyben ábrázoljunk a segítségével (6. rész, 27. oldal), és végül megvizsgáljuk annak a lehetőségeit, hogy hogyan lehet az így kialakított struktúrát a lényeges igei szerkezetek felismerésére felhasználni (7. rész, 34. oldal).

## 2 Hozzávalók – absztrakt tárgyalás

Képzeljünk el egy komódot, aminek három fiókja van. A piros színű fiókban egy könyv van, a fehér fiókban egy labda, a zöld fiók pedig üres. Van tehát egy fő elem (a komód), tartoznak hozzá egyértelműen azonosítható helyek (a három különböző színű fiók), és vannak dolgok (a könyv és a labda), amik ezeket a helyeket elfoglalhatják. Egy helyen maximum egy dolog lehet, az is előfordulhat, hogy valamelyik hely üres. Ez az egyszerű szemléletes bemutatása annak az alapmodellnek, amit a tanulmányban végig használni fogunk. Nézzük ezt meg alaposabban formálisan.

### 2.1 Alapelemek

A következő terminusokat fogjuk használni. A fő elem elnevezése **forrás**, **gyökér** vagy **minimum**. A forráshoz kapcsolódó helyeket egyszerűen **helynek** nevezzük, a dolgokat pedig, melyek a helyeken megjelenhetnek, **kitöltőnek**. A forrást általában  $i$ -vel, a helyeket az ábécé elejéről vett kisbetűkkel ( $e, f, \dots$ ), a kitöltőket az ábécé végéről vett kisbetűkkel ( $w, x, \dots$ ) jelöljük. E háromféle entitásra – forrás, hely és kitöltő – együttesen **elemekként** utalunk. Azokat a helyeket, melyeken nincs kitöltő **szabad helyeknek**, azokat pedig, melyeken van kitöltő **kitöltött helyeknek** nevezzük. A példabeli komódban egy szabad és két kitöltött hely van.

### 2.2 A két művelet

Definiálunk két műveletet, melyek a különböző típusú elemek közötti lehetséges kapcsolódási pontokat írják le. Mindkét műveletnek lesz egy kifejtettebb, műveleti jelekkel operáló megjelenítési módja, és egy másik, tömörebb, ahol a térbeli elrendezés ábrázolja a művelet eredményét és a résztvevő elemek közötti viszonyokat. A műveletsorokat balról jobbra olvassuk és hajtjuk végre.

Az első a **hely-hozzáadás**, mikor a forráshoz egy helyet kapcsolunk hozzá – azaz készítünk a komódunkhoz egy újabb fiókot.  $i + e$  a jele annak, hogy az  $i$  forráshoz (vagy mint látni fogjuk tetszőleges szerkezethez) hozzáadjuk az  $e$  helyet. Egy forráshoz tetszőleges számú helyet hozzá lehet adni. Az  $ie$  térbeli elrendezés (egyszerűen az egymás mellé írás) mutatja, hogy az  $i$  forráshoz (ez szerepel mindig elől) hozzá van kapcsolva egy  $e$  hely.

A másik művelet a **hely-kitöltés**: ekkor egy meglévő helyhez egy kitöltőt rendelünk hozzá, egy helyre „beteszünk” egy kitöltőt (pl.: cipőt a kék fiókba).  $e \curvearrowright w$  a jele annak, hogy az  $e$  helyet kitöltjük a  $w$  kitöltővel. Egy helyhez csak egy kitöltőt lehet ezen a módon hozzárendelni. Ezt a műveletet a  $e$  térbeli elrendezéssel (fent a hely, alatta a kitöltő) is ábrázolni fogjuk.

Egyik műveletnél sem teszünk különbséget a műveleti eljárás („kitöltjük”) és a művelet eredményeként létrejövő viszony („ki van töltve”) között. Szükség szerint használ(hat)juk a műveletek procedurális és deklaratív értelmezését.

### 2.3 A szerkezet definíciója

Az egy forrásból,  $e$  forráshoz hozzáadott 0 vagy több helyből, és 0 vagy több meghatározott helyhez hozzárendelt kitöltőből álló összetett entitásokat **szerkezetek**nek nevezzük. A következő példa egy 2 kitöltött hellyel bíró szerkezetet mutat be:

$$i + e \frown w + f \frown x = \underset{wx}{ief}$$

A műveletek értelemszerűen vonatkoztathatók szerkezetekre is:  $+$  a szerkezet forrására,  $\frown$  pedig a szerkezet egy adott helyére értendő. Ha a szóban forgó szerkezetnek több szabad helye van (pl.:  $ief$ ), akkor a művelethez írt alsó indexszel lehet egyértelművé tenni, hogy a kitöltés melyik helyre vonatkozik:

$$ief_z = ief \frown_e z \neq ief \frown_f z = \underset{z}{ief}$$

Azonos helyre vonatkozó két művelet értelemszerűen nem cserélhető fel, kizárólag (1) hely-hozzáadás, (2) hely-kitöltés sorrendben végezhető el, a két különböző helyre vonatkozó műveletek viszont felcserélhetők, elvégzésük sorrendje tetszőleges, például:

$$i + e \frown_e w + f = i + e + f \frown_e w = \underset{w}{ief}$$

Megjegyezzük, hogy a fent definiált két művelet nem felel meg a művelet szűken vett klasszikus matematikai fogalmának, mert különféle típusú entitásokkal dolgozik. Mindkettő ún. külső művelet, mely az egyik argumentumát külső halmazból – esetünkben nem a szerkezetek közül – veszi. A hely-hozzáadás egy szerkezetből és egy helyből készít egy újabb szerkezetet, az (adott helyre vonatkozó) hely-kitöltés pedig egy szerkezetből és egy kitöltőből készít egy újabb szerkezetet.

### 2.4 Fogalmak, függvények, aszimmetria

Azt a szerkezetet, melyben minden hely ki van töltve **nyelőnek**, illetve a későbbiekben **maximumnak**, **maximális szerkezetnek** is fogjuk nevezni.

Néhány, szerkezeteken értelmezett, nemnegatív egész értékű, egyszerű függvényt vezetünk be. A **helyek száma** ( $h$ ) függvény megadja, hogy hány hely van a szerkezetben; a **kitöltöttség** ( $k$ ) függvény megadja, hogy hány helyen van kitöltő a szerkezetben; a **hossz** ( $l$ ) függvény pedig azt, hogy összesen hány elem van a szerkezetben.

Mindig érvényes az ún. aszimmetria tulajdonság:  $k \leq h$ , vagyis hogy a hely-kitöltések száma legfeljebb annyi lehet, mint a hely-hozzáadások száma ( $|\frown| \leq |+|$ ), mivel csak meglévő helyet lehet kitölteni. Ezenkívül az elem definíciójából adódóan nyilván mindig érvényes a fenti függvényeknek az alábbi egyszerű összefüggése:  $l = h + k + 1$ , ahol a  $+1$  a forrás miatt szerepel.

$h$  és  $k$  aszimmetria tulajdonságból adódó speciális viszonyaira az alábbi megnevezéseket használjuk: **teljesen kitöltött (tk)** egy szerkezet, ha  $k = h$ ; **szabad hellyel bíró (szhb)**, ha  $k < h$ . E két halmaz diszjunkt, és uniójuk az összes szerkezet. Az utóbbi halmazon belül

$k$	$h = 0$	1	2	3
0	tk	$\frac{\text{tklen}}{\text{szhb}}$	$\frac{\text{tklen}}{\text{szhb}}$	$\frac{\text{tklen}}{\text{szhb}}$
1	×	tk	$\frac{\text{vgys}}{\text{szhb}}$	$\frac{\text{vgys}}{\text{szhb}}$
2	×	×	tk	$\frac{\text{vgys}}{\text{szhb}}$
3	×	×	×	tk

1. táblázat.  $h$  és  $k$  viszonyából adódó szerkezetípusok: teljesen kitöltött (tk); szabad helyvel bíró (szhb); teljesen kitöltetlen (tklen); vegyes.

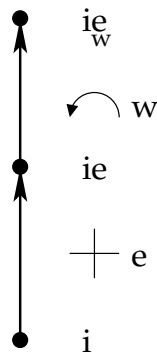
elkülöníthetjük a **teljesen kitöltetlen (tklen)** szerkezeteket, melyekre  $k = 0$  és  $h > 0$ , illetve a **vegyes** szerkezeteket, melyekre  $0 < k < h$ , azaz szabad és kitöltött helyük is van (1. táblázat).

Azt mondjuk, hogy  $a$  szerkezet **illeszkedik**  $b$ -re, ha rövidebb, azaz  $l(a) \leq l(b)$ , és  $a$  helyei és az adott helyeken lévő kitöltők  $b$ -ben is megvannak.  $i$  illeszkedik  $ie$ -re,  $ie$  illeszkedik  $ie_w$ -re,  $ie$  illeszkedik  $ief$ -re, de nem illeszkedik  $ie$ -re vagy  $if$ -re. Ha  $a$  illeszkedik  $b$ -re, akkor azt mondjuk, hogy  $a$   $b$ -nek **alszerkezete**. Látjuk, hogy a fentiek alapján a maximális és a tk szerkezet fogalma azonos. A tk megjelölést inkább adott maximális szerkezet tk alszerkezeteire fogjuk használni.

Bevezetjük a **valódi szerkezetek**nek a jelen absztrakt tárgyalás során még meglehetősen üres fogalmát is. Ez egyelőre csupán annyit jelent, hogy egy maximális szerkezet alszerkezetei közül bizonyos szerkezeteket valódiként külön megjelölünk.

## 2.5 Hálók

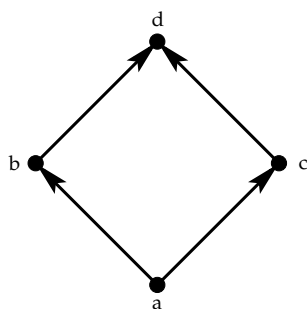
A teljesen kitöltött, maximális szerkezeteket ábrázolhatjuk irányított gráfokként: a kiinduló csomópont a forrás lesz; az irányított éleken a korábban bevezetett két szerkezetépítő művelet valamelyike fordulhat elő; a további csomópontokban pedig mindig az adott művelet eredményeként létrejövő szerkezet szerepel (1. ábra). Ilyen gráfok esetében a csomópont és a szerkezet kifejezéseket szinonimaként használhatjuk.



1. ábra. Az egy hellyel rendelkező teljesen kitöltött  $i + e \curvearrowright w = ie_w$  szerkezet ábrázolása gráfként. A  $+$  és  $\curvearrowright$  művelet meghatározását ld. a 2.2. részben. Erre a struktúrára fogunk még hivatkozni a továbbiakban, és az  $A$  jelölést fogjuk használni rá.

A nyilakat mindig fölfelé fogjuk irányítani, a forrás így a gráfára alján, a nyelő pedig a tetején fog megjelenni. Vegyük észre, hogy a gráf minden csomópontja szerkezet, mégpedig a maximális szerkezet alszerkezete. A nyilak mindig az eggyel hosszabb szerkezet felé mutatnak.

Felidézzük az algebrából ismert **háló** fogalmat. Mivel csak véges hálóink lesznek, csak azokkal foglalkozunk. Két egymással ekvivalens definíció létezik, egy olyan, ami a részbenrendezés fogalmán alapul, és egy olyan, ami szerint a háló egyfajta algebrai struktúra. Az első definíció szerint a háló egyrészt egy részbenrendezett halmaz, másrészt pedig van legkisebb eleme (minimuma) és legnagyobb eleme (maximuma). Ez annyit jelent, hogy van rajta egy  $\leq$  reláció, ami olyan mint egy rendezés (reflexív, antiszimmetrikus és tranzitív), de nem feltétlenül hasonlítható össze általa minden elempár; a kitüntetett legkisebb (illetve legnagyobb) elemmel viszont minden elem összehasonlítható, és mind nagyobb (illetve kisebb) nála (2. ábra).



2. ábra. Egy egyszerű háló. A szokásnak megfelelően nyilak ábrázolják a részbenrendezési relációt, a nyilak a nagyobb elem felé (és mindig fölfelé) irányulnak. Így igaz, hogy  $a \leq b$ ,  $b \leq d$ ,  $a \leq c$ ,  $c \leq d$ , valamint a tranzitivitás miatt  $a \leq d$ , a  $b$  és a  $c$  azonban nem összehasonlítható. A minimum  $a$ , a maximum  $d$ .

Ezen első definíció matematikailag precíz formája, illetve a második definíció, mely szerint a háló, egy két darab kétváltozós művelettel – minimumképzés ( $\wedge$ , meet) és maximumképzés ( $\vee$ , join) – bíró algebrai struktúraként is megragadható, megtalálható például Partee et al. 1990: 11.2 részében. Szemléletesen: a háló egy irányítottkör-mentes gráf, amiben a nyilak fölfelé mutatnak, és van egyértelmű – egy pontként megjelenő – „alja” és „teteje”.

Vegyük észre, hogy a teljesen kitöltött szerkezeteket ábrázoló gráfok (1. ábra) hálók: nyilván nincs bennük irányított kör (emiatt érvényes a részbenrendezettség), és van legkisebb elem (a forrás) és legnagyobb elem (a nyelő). Most már világos, hogy az egyes csomópontok elnevezését eleve a vonatkozó hálóelméleti fogalomnak megfelelően választottuk. Az említett két hálóművelet ( $\wedge$ ,  $\vee$ ) nem keverendő össze a 2.2. részben bevezetett két szerkezetépítő művelettel ( $+$ ,  $\curvearrowright$ ). Utóbbiakkal „fölfelé” haladunk a hálókön, azaz mindkét szerkezetépítő művelet a maximumképzés ( $\vee$ ) speciális esetének tekinthető. A hálók szokásos értelemben vett hossza („a leghosszabb irányított út csomópontjainak száma”) megegyezik a maximális szerkezet  $l$  szerinti hosszával. A szerkezetekre fent bevezetett illeszkedés fogalom mindössze annyit jelent, hogy összehasonlítható a két elem, azaz  $a$  illeszkedik  $b$ -re pontosan akkor, ha  $a \leq b$ .

A hálóink **szintezhetők**, azaz meg tudunk adni egy olyan ( $r$ ) függvényt, mely megfelel a rendezésnek, azaz:  $a \leq b$  és  $a \neq b$  esetén  $r(a) < r(b)$ . A már megismert  $l$  függvény éppen jó erre a célra. A hálók felrajzolhatók úgy, hogy az azonos szinten lévő (azaz azonos  $r$  értékű) csomópontok egy vízszintes egyenesre essenek. A hálókat a legtöbb esetben így fogjuk felrajzolni.

### 3 Példák szerkezethálókra

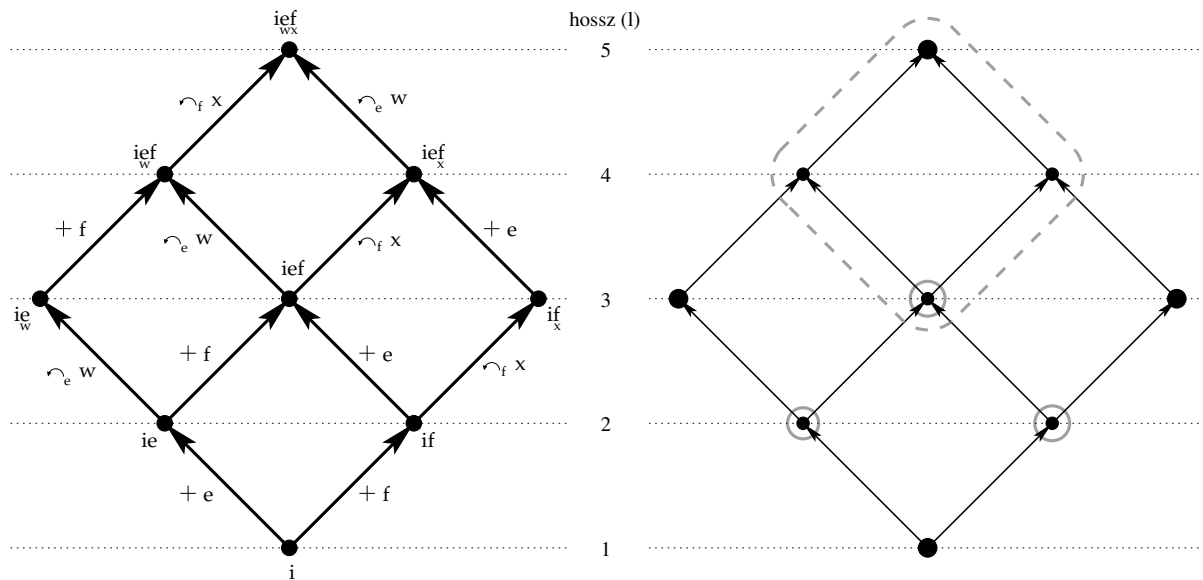
Ebben a részben bemutatjuk az 1, 2 és 3-helyes ( $h = 1, 2, 3$ ) teljesen kitöltött szerkezetet ábrázoló hálót, és megvizsgáljuk tulajdonságaikat.

#### 3.1 Az 1-helyes háló

Az 1-helyes ( $h = 1$ ) háló azonos az 1. ábrán (15. oldal) szereplő gráffal. Ahogy az ábránál is utaltunk rá, ez lesz a szerkezetháló között az alapegységünk, és  $A$ -ként fogunk hivatkozni rá. Három csomópontot, egy 1, egy 2 és egy 3 hosszúságú szerkezetet tartalmaz, a csomópontok száma szintenként ( $r$  szerint emelkedően): 1 1 1.

#### 3.2 A 2-helyes háló

Bonyolultabb és egyúttal érdekesebb szerkezetű a 3. ábrán látható 2-helyes háló.



3. ábra. (balra) A 2-helyes háló a csomópontok (szerkezetek) és élek (műveletek) formális megjelenítésével. (jobbra) A 2-helyes háló különféle tulajdonságú csomópontjainak bemutatása: (1) nagy korong = teljesen kitöltött; (2) kis korong = szabad helyet bíró; (3) bekarikázott kis korong = teljesen kitöltetlen; (4) bekarikázatlan kis korong = vegyes; (5) szaggatott vonallal körbevett rész = nemtriviális szerkezetek (meghatározását ld. alább a szövegben).

Ez a háló egy 2-helyes maximális szerkezet ( $ief$ ) alszerkezeteit írja le. A hálóban lévő 4 hosszúságú utak adják meg azokat a műveletsorrendeket, melyek által a forrástól a maximális szerkezetig juthatunk el. A 3. ábra bal alsó „éle” pontosan megfelel a  $h = 1$  esetnek (azaz az 1. ábrának). Látjuk, hogy  $h = 1$  esetén csak teljesen kitöltött és teljesen kitöltetlen csomópontok fordulnak elő, jelen ábrán a többi csomóponttípus is megjelenik.

A  $2 \times 2$ -es négyzetháló struktúra pontosan a két szerkezetépítő művelet felcserélhetőségi viszonyainak (ld. 14. oldal) következményeképpen áll elő, abból adódik, hogy az adott helyre vonatkozó hely-hozzáadás mindenképpen megelőzi a hely-kitöltést. Ennek megfelelően van az egyes csomópontokban egy vagy két kifelé mutató nyíl. Formálisan és a címkézéstől eltekintve az 1. ábráról ismert  $A$  háló önmagával való direkt szorzatát, azaz direkt második hatványát

( $A \times A = A^2$ ) látjuk a 3. ábrán. A csomópontok száma 9 ( $3^2$ ), a csomópontok száma szintenként ( $r$  szerint emelkedően): 1 2 3 2 1.

A későbbiekben szükség lesz az adott szerkezetenél „1-gyel kisebb” szerkezetek meghatározására. Természetesen a háló részbenrendezése ( $\preceq$ ) szerinti 1-gyel kisebbségről van itt szó (vö: 2. ábra, 16. oldal). Ezek nyilván azok a csomópontok, ahonnan az adott csomópontba nyíl mutat. Most csak annyit jegyzünk meg, hogy sok esetben ez nem egyértelmű, több ilyen szerkezet is lehet: egész pontosan  $h$  darab, azaz amennyi hely van az adott szerkezetben. Ugyanis adott szerkezet minden helyén mindig pontosan egy szerkezetépítő művelet inverzét végezhethetjük el: ha ki van töltve a hely, akkor szabaddá tesszük, ha nincs, akkor pedig elhagyjuk a helyet. Azokat a szerkezeteket, melyekhez több „1-gyel kisebb” szerkezet létezik, **nemtriviális** szerkezeteknek nevezzük, ez pontosan  $h \geq 2$  esetén teljesül (ld. a 3. ábrán a szaggatott vonallal körbevett részt).

### 3.3 A 3-helyes háló

A 3-helyes háló kényelmesen csak 3 dimenzióban ábrázolható. A 4. ábrán látható rajzot (Epstein 1993) 104. oldaláról vettük át, de korábban megjelenik például (Epstein 1960) 309. oldalán is.

A struktúra minden dimenzió mentén 2-2 (3-dimenziós) kocka egymásra építésével áll elő, ennek megfelelően szemléletesen (3-dimenziós) **duplakockának** is fogjuk nevezni. (Hasonlóan: az 1-helyes és a 2-helyes hálót 1- illetve 2-dimenziós duplakockának is nevezzük, mivel pontosan ugyanilyen módon épül fel 1- és 2-dimenziós kockákból.) A háló tényleges kockákkal való ábrázolásának előnyös tulajdonsága, hogy (1) térben nincs élkereszteződés; (2) ténylegesen egy magasságban vannak a háló azonos szinten lévő (azonos  $r$  értékű) csomópontjai, értsd: fizikailag valóban lehetséges úgy tartani egy ilyen térbeli alakzatot, hogy egy magasságban legyenek. Az ábrán egymáshoz közel rajzolt pontok – például  $(0, 2, 0)$  és  $(1, 0, 1)$  – egy szinten vannak.

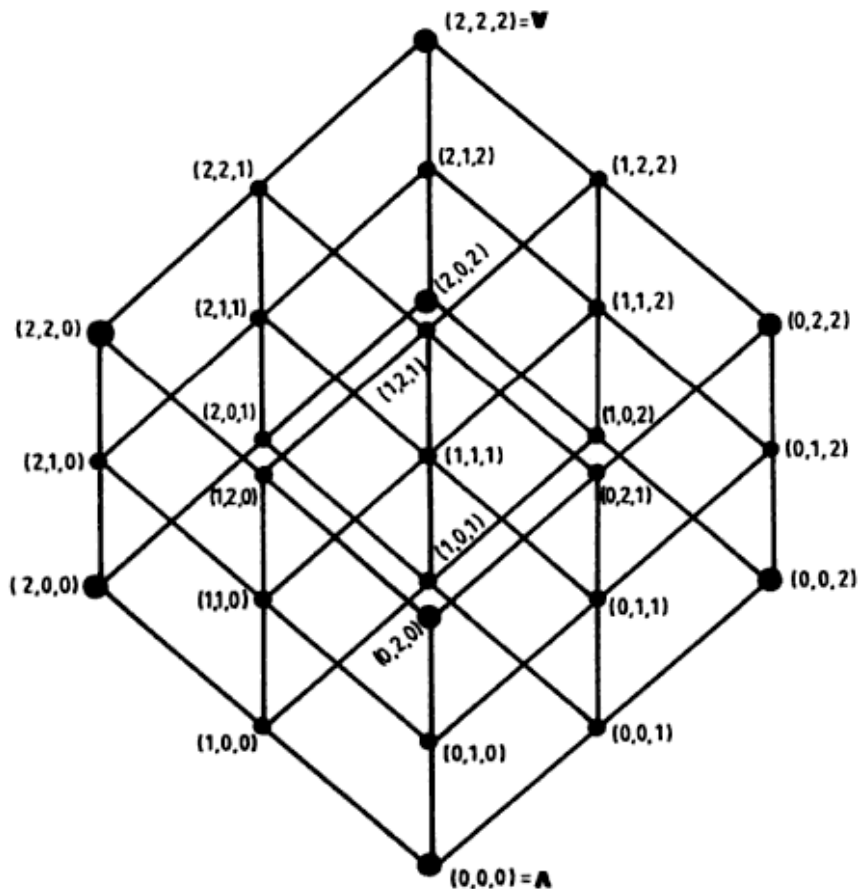
Említettük (17. oldal), hogy az 1-helyes  $A$  háló önmagával való direkt szorzataként ( $A^2$ ) kapjuk meg a 2-helyes hálót. Most azt látjuk, hogy a 3-helyes háló pedig éppen  $A \times A \times A = A^3$ . A csomópontok száma 27 ( $3^3$ ), a csomópontok száma szintenként ( $r$  szerint emelkedően): 1 3 6 7 6 3 1.

### 3.4 A duplakockák általános tulajdonságai

Vegyük észre, hogy ezek a hálók egy sorozatot alkotnak: a  $h$  helyes háló, azaz a  $h$  dimenziós duplakocka, mindig  $A$ -nak  $h$ -adik direkt hatványa lesz. Az  $A^h$  hálók tulajdonságait a 2. táblázatban látjuk.

A táblázatban gyorsnak neveztük a  $3^n$  szerint (vagy még gyorsabban) haladó sorozatokat, *lassúnak* a  $2^n$  szerint haladókat. A teljesen kitöltött csomópontok fontos szerepet fognak még játszani (vö: 6.4. rész), látjuk, hogy ezek lassú sorozatot alkotnak, azaz „viszonylag kevés” van belőlük. A táblázatbeli sorozatokhoz képest még sokkal lassabb az  $A^h$ -ban az adott csomópont-hoz tartozó „1-gyel kisebb” csomópontok maximális száma: konkrétan  $h$ .

Láttuk, hogy a  $A^h$  azonos szintű csomópontjainak elemszáma  $h = 1$  esetén: 1 1 1;  $h = 2$  esetén: 1 2 3 2 1;  $h = 3$  esetén: 1 3 6 7 6 3 1. Ezek a számcsoportok éppen az ún. trinomiális háromszög sorai (ld. az A027907 számú sorozatot az OEIS-ben (OEIS Foundation Inc. 2011)). Ennek a háromszögnek az egyes értékei úgy kaphatók meg, hogy az adott pozíció fölött elhelyezkedő 3 számot összeadjuk. Ez a hármas („tri”) elrendeződés a két szerkezetépítő művelet



4. ábra. A 3-helyes háló (az ábra forrása: (Epstein 1993: 104)). Ez egyszerűen egy csúcsára állított 8 kiskockából álló  $2 \times 2 \times 2$ -es kocka síkban megjelenített rajza. Térben elképzelve a  $(0, 2, 0)$  jelű pont fekszik hozzánk legközelebb, a  $(2, 0, 2)$  tőlünk legtávolabb. Az ábrán szereplő kódok könnyen megfeleltethetők a korábban bevezetett jelöléseknek: az egyes kódpozíciók megmutatják rendre az  $e, f, g$  helyekről, hogy a hely: nincs (0); van és szabad (1); vagy van és ki van töltve (2). A  $(2, 1, 0)$  például a  $ief$ -nek felel meg. A szerkezetek hossza ( $l$ ) éppen a kódpozíciók összege +1.

tulajdonságaiból, hierarchiájából adódik, ugyanis minden hely tekintetében 3 lehetőség van: nincs a hely; van és szabad; van és ki van töltve. Így aztán  $h$  hely esetén éppen  $3^h$  lehetőség van, ennyi az  $A^h$  csomópontjainak száma, és természetesen a trinomiális háromszög sorösszege is.

Minden  $h$ -dimenziós duplakocka elhelyezhető úgy a  $h$ -dimenziós térben, hogy az azonos szinten lévő (azonos  $r$  értékű) csomópontjai egy magasságban legyenek, a forrás a duplakocka legalulra eső csúcsán, a maximális ( $h$  helyes tk) szerkezet pedig a legfelülre eső csúcson legyen. A teljesen kitöltött csomópontok bármennyi dimenzióban a duplakocka csúcsain helyezkednek el, a  $h$  helyes tklen szerkezet pedig a duplakocka „középpontjára” fog esni. Megfigyelhető, hogy  $h \leq 3$  esetén ez a középpont a  $h$ -dimenziós duplakockának az *egyetlen* nemfelszíni pontja; sejtés, hogy ez  $h > 3$  esetén is így van.

Hogyan lehet megjeleníteni, lerajzolni  $h > 3$  esetén az  $A^h$  hálókat? Nincs egyszerű megoldás. A „triviális” ábrázolási mód a  $2 \times 2 \times \dots \times 2 = 2^h$  darab  $h$ -dimenziós kiskockából álló



$A^h \dots$	formula	$h = 0$	1	2	3	4	5	OEIS azon	megj
csomópontjai	$3^h$	1	3	9	27	81	243	A000244	gyors
élei	$2h \cdot 3^{h-1}$	0	2	12	54	216	810	A212697	gyors
tk csp-jai	$2^h$	1	2	4	8	16	32	A000079	lassú
szhb csp-jai	$3^h - 2^h$	0	1	5	19	65	211	A001047	gyors
tklen csp-jai	$2^h - 1$	0	1	3	7	15	31	A000225	lassú
vegyes csp-jai	$3^h - 2^{h+1} + 1$	0	0	2	12	50	180	A028243	gyors
nemtriviális csp-jai	$3^h - 2h - 1$	0	0	4	20	72	232	A061981	gyors
azonos szintű csp-jai (trinomiális háromszög)	$T(h, n) = \sum_{i=0,1,2} T(h-1, n-i)$				1	3	6	...	A027907
			1	2	6				
			1	3	7				
			1	2	6				
				1	3				
					1				

2. táblázat. Az  $A^h$  hálók egyes jellemzői  $h$  függvényében. Az utolsó előtti oszlopban az adott számsorozat OEIS-beli (OEIS Foundation Inc. 2011) azonosítója szerepel.

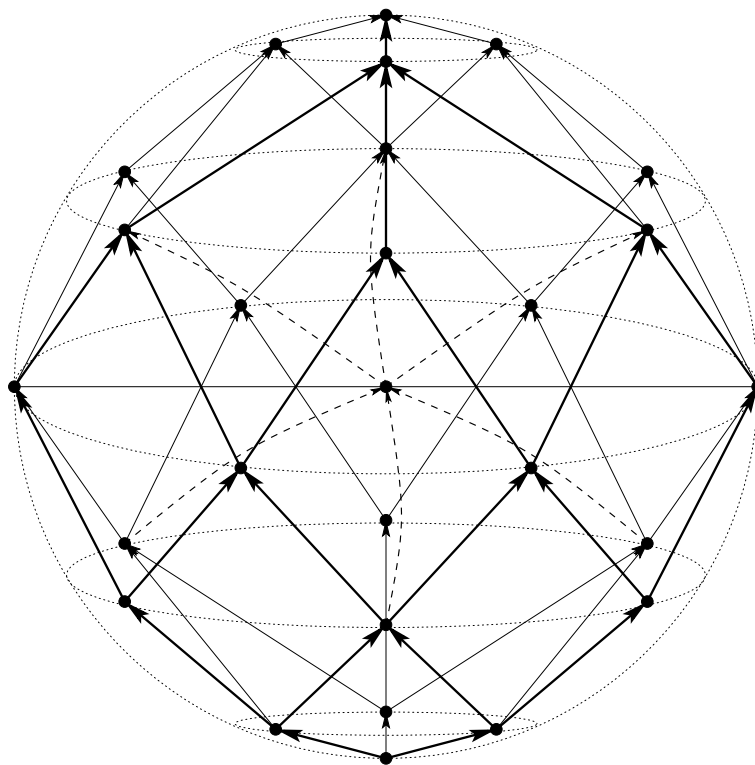
$h$ -dimenziós duplakocka, melynek kényelmes ábrázolásához  $h > 3$  dimenzió szükséges. Ha igaz az, hogy egyetlen nemfelszíni pont van, akkor elegendőnek tűnik a középpont plusz egy gömbfelület, azaz gömbrerajzolhatók a hálók, és ezzel egy általános, minden  $h$ -ra működő alternatív ábrázolási módszert kapunk. Ez a módszer  $h \leq 3$ -ra valóban működik is (5. ábra), nagyobb dimenzióra azonban a gyakorlatban ez sem megoldás, mert az élkereszteződés nélküli ábrázoláshoz még mindig szükség van egy  $h$ -dimenziós gömbfelületre.

#### 4 Post-hálók és Boole-hálók

Epstein (1993: 104) a tárgyalt hálókat (1., 3., és 4. ábra) **harmadrendű Post-háló**nak (*Post lattice of order 3*) nevezi. E hálókat a Boole-háló általánosításaként vezeti be.

A Boole-háló a két pontból álló háló direkt hatványaiként állnak elő, szemléletesen „szimplakockaként” jeleníthetők meg: a  $h = 2$ -re négyzetet kapunk (2. ábra), a  $h = 3$ -ra kockát, a  $h = 4$ -ra tesszeraktot (azaz négydimenziós kockát). A Boole-háló felfogható egy  $h$  elemű halmazon értelmezhető kétértékű (például logikai) függvények rendszereként. Ilyen kétértékű függvény például a részhalmazképzés (a kiválasztott elemekhez 1-et/igazat rendelünk, a többihez 0-t/hamist), ezáltal a Boole-háló felfogható egy  $h$  elemű halmaz részhalmazáiból képzett hálóként is, a klasszikus metszet ( $\cap$ ) és unió ( $\cup$ ) műveletekkel.

Az általánosítási lépés talán úgy ragadható meg legkönnyebben, hogy a kiinduló elem nem kételemű, hanem  $n$ -elemű lánc, az  $n$ -edrendű Post-háló az  $n$ -elemű lánc direkt hatványai. A Post-háló véges, teljes, disztributív (Epstein 1993: 102), szintezhető háló. Nyilvánvaló, hogy a Boole-háló azonosak a másodrendű Post-hálók (Pagliani és Chakraborty 2008: 254). Ahol a klasszikus kétértékű logikában Boole-hálókat használunk, ott az  $n$ -értékű logikában  $n$ -edrendű Post-hálókat (Epstein 1993 és Chechik et al. 2001: 1/b. ábra). A részhalmazképzés 3 értékre való általánosítását elképzelhetjük úgy, hogy van egy harmadik, „félíg tartozik bele”



5. ábra. A 3-helyes háló, azaz a 3-dimenziós duplakocka ábrázolása gömbre rajzolva.

érték is. Valami hasonlót látunk a duplakockák esetében is: amint már többször volt róla szó, a  $+$  és  $\wedge$  művelet tulajdonságainak köszönhetően minden hely tekintetében háromféle állapot fordulhat elő. Világos tehát, hogy éppen a harmadrendű Post-hálókkal van dolgunk: a duplakockák azonosak (izomorfak) a harmadrendű Post-hálókkal (avagy 3-Post-hálókkal). Amiatt, hogy a Boole-háló az  $n = 2$  eset, a 3-Post-háló pedig az  $n = 3$  eset, nevezik az utóbbit 'Boolean trilattice'-nek (kb. Boole-hármasháló) is (Biedermann 1998). (A terminológia itt nem egységes. Előfordul az is, hogy a 'trilattice' terminust a  $h = 3$  esetre használják (Shramko et al. 2001).)

A továbbiakban a duplakockákra a ( $h$  darab helyre) vonatkoztatott harmadrendű Post-háló, röviden **3-Post-háló  $h$ -ra** vagy  **$h$ -dimenziós 3-Post-háló** hivatalos megjelölést is használhatjuk. Eddigi példák: 3-Post-háló 1-re (1. ábra), 3-Post-háló 2-re (3. ábra), 3-Post-háló 3-ra (4. ábra). Értelmezzük szemléletesen az utóbbi ábra esetén a két hármast. Az első hármast azt adja meg, hogy egy „duplakockaélen” hány csomópont (azaz hányféle „érték”) van, ez a Post-háló rendje, ami esetünkben fixen 3; a második hármast azt adja meg, hogy a forrásból (avagy gyökérből) hány nyíl indul ki, azaz hogy hány helyre vonatkozik a Post-háló, ez változik, itt épp  $h = 3$ .

Szemléletesen a Boole-háló  $h$  darab kétállású kapcsoló lehetséges állapotait (és a közöttük egy kapcsolással megvalósítható állapotváltozásokat) jeleníti meg, a 3-Post-háló pedig ugyanezt teszi, csak háromállású kapcsolókkal. Így akár hívhatjuk **uszodaihajszerítés-háló**-nak is, magunk elé képzelve az egy helységben működő  $h$  darab gyenge és erős fokozattal is rendelkező hajszerítőt.

És ez végülis ugyanaz, mint a tanulmány elején említett komód-hasonlat.

## 5 Igei szerkezetek és hálók

A bevezetőben azt ígértük, hogy az igei szerkezeteknek egy új formális modelljét fogjuk bemutatni. Egész eddig egy absztrakt modellel foglalkoztunk, igei szerkezetekről alig esett szó. Most tesszük meg azt a fontos lépést, hogy az igei szerkezeteket az imént felépített absztrakt szerkezetmodell segítségével ábrázoljuk. Vajon hogy jönnek a képbe a tanulmány címében már elővételezett igei szerkezetek és egyáltalán az igeik? Erről szól ez a fejezet.

### 5.1 Megfeleltetés

Megadunk egy megfeleltetést a modell 2. részben ismertetett alkotórészeinek vonatkozásában a következők szerint. A gyökér egy ige lesz; a helyek esetragok (vagy névutók); a kitöltők pedig főnevek (névszók). Egész pontosan a helyek az ige bővítményeinek esetragjai (vagy névutói), a kitöltők pedig az ige bővítményeiként megjelenő főnevek (névszók) lesznek.

A szerkezeteket a gyökérről nevezzük el, tehát az így konkretizált szerkezeteket **igei szerkezeteknek** fogjuk nevezni. A fenti megfeleltetéssel kapunk az absztrakt szerkezetekből igei szerkezeteket. A korábbiaknak megfelelően az igét általában  $i$ -vel, az esetragokat/névutókat  $e, f, \dots$  kisbetűkkel, a kitöltő névszókat pedig  $w, x, \dots$  kisbetűkkel jelöljük. Az igei szerkezetek alapelemei abban az értelemben egységes halmazt képeznek, hogy minden alapelemről elmondhatjuk, hogy egy plusz információegységet ad hozzá a szerkezethez, legyen az nyelvi szempontból tartalmi vagy strukturális.

Az alapelemek fenti egyszerű megfeleltetésének, illetve annak köszönhetően, hogy az igei szerkezetek az absztrakt szerkezeteknek egy konkrét megvalósulását, speciális esetét jelentik, minden, amit a korábbi fejezetekben az absztrakt szerkezetekről mondtunk, érvényes a most bevezetett igei szerkezetekre is. Az alapelemek ismeretében minden kijön az absztrakt modellből, mindent készen kapunk, a két szerkezetépítő művelettel, a leírt fogalmakkal, aszimmetriával, függvényekkel, illeszkedéssel és a harmadrendű Post-hálók összes tulajdonságával együtt.

Nézzük egy példát. Vegyünk egy egy tagmondatból álló egyszerű mondatot.

*Lencsi könyvet olvas.*

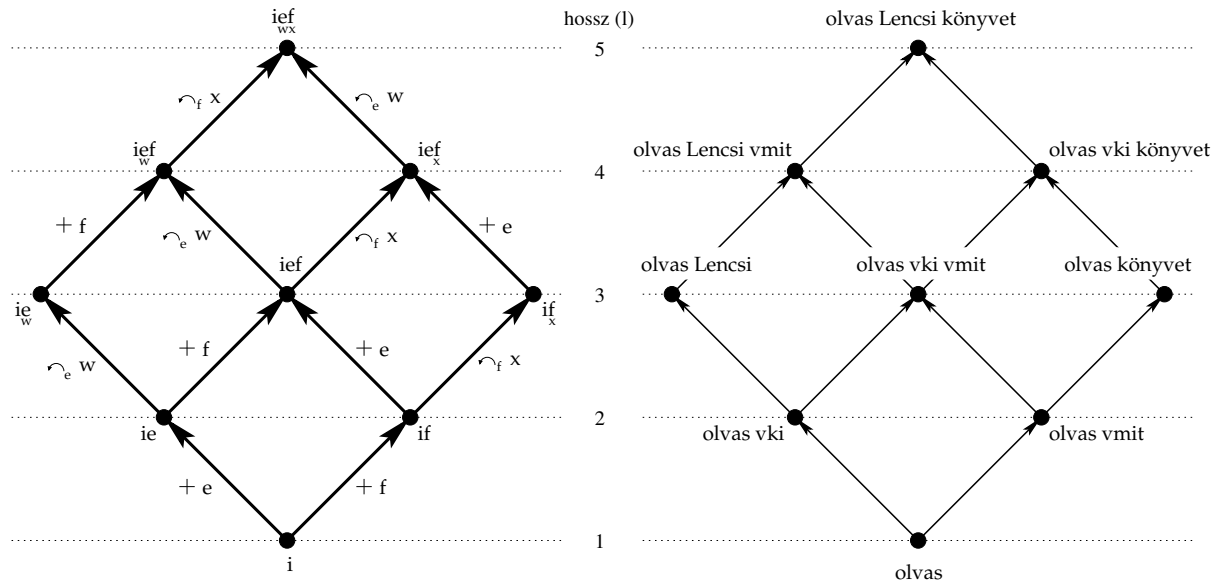
Van egy igénk: *olvas* ( $i$ ); két hely az ige mellett: egy alanyesettel megjelölt hely ( $e$ ) és egy tárgyesettel megjelölt hely ( $f$ ); valamint két kitöltő: az  $e$  helyen lévő *Lencsi* ( $w$ ) és az  $f$  helyen lévő *könyv* ( $x$ ). Azaz ez az egyszerű mondat formálisan a korábban (2.3. és 3.2. rész) bemutatott klasszikus két hellyel bíró szerkezetként ragadható meg:  $i e f$ . A szerkezetet előállító két helyhozzáadást és két hely-kitöltést tartalmazó  $i + e \frown w + f \frown x$  művelet sor így értelmezhető: veszünk egy igét, hozzáteszünk egy alanyi helyet, erre a helyre berakunk egy konkrét alanyt, hozzáteszünk az igehez egy tárgyi helyet is, és oda berakunk egy tárgyat. Alapesetben egy igei szerkezet alapelemei mindig egy tagmondaton belül helyezkednek el.

Az egyszerű mondat központi eleme az ige, az egyszerű mondat leegyszerűsítve egy igéből és az ige bővítményeiből áll, az egyszerű mondat az igének és az ige bővítményeinek összessége. Lényegében ennek felel meg, ezt ragadja meg a fenti formális leírás – a fentiek szerint első körben csak a névszói csoport és a névutós csoport bővítményeket tekintetbe véve. Ilyen bővítményből persze lehet a tagmondatban 4, 5 vagy akár több is, ezért érintettük a  $h > 3$ -ra vett („magasabb dimenziós”) 3-Post-hálókat az absztrakt tárgyalás során (3.4. rész). A továbbiakban mikor mondatról beszélünk, egyszerű mondatot vagy tagmondatot értünk alatta.

A maximális szerkezetet az igei szerkezetek esetében **mondatváz**nak fogjuk nevezni. Azt hangsúlyozandó, hogy nem konkrét (tag)mondatokról van szó a mondatvázakat általában *igetű*

+ helyek sorrendben fogjuk feltüntetni. A fenti *Lencsi könyvet olvas.* mondat mondatváza tehát: *olvas Lencsi könyvet* (mondatvégi írásjel nélkül). Valójában ezt, a mondatvázat, írja le a fenti igei szerkezet, nem a konkrét mondatot. A mondatváz nem csak sorrendi konvenció, hanem valódi általánosítás a mondatok felett, hiszen számos konkrét mondatnak ugyanaz a mondatváza: az *Egy könyvet olvas Lencsi.* és a *Lencsi olvasta a könyveket.* mondatnak például szintén a fenti.

Az eddigiek szemléltetésére lássuk a 6. ábrát.



6. ábra. A 2-helyes mondatvázat (és alszerkezeteit) bemutató háló. (balra) A csomópontok és élek formalizált megjelenítése a 3. ábra bal oldaláról átvéve. (jobbra) A bal oldalon lévő absztrakt modellt megvalósító *olvas Lencsi könyvet* 2-helyes mondatváz hálójá:  $i = \text{olvas}$ ;  $e = \emptyset$  (alany);  $f = -t$  (tárgy);  $w = \text{Lencsi}$ ;  $x = \text{könyv}$ . A *vki/vmi* megjelölés itt és a továbbiakban is a kitöltetlen helyet jelenti. A háló csomópontjai igei szerkezetek, a maximális igei szerkezet alszerkezetei. A minimális igei szerkezet maga az ige, belőle mint gyökérből „virágzanak ki” a különféle igei szerkezetek. Az ábra természetesen egy 2-dimenziós harmadrendű Post-háló.

Az 1-helyes *olvas Lencsi* mondatvázat a 6. ábrán látható háló bal alsó „élét” alkotó három csomópontból álló háló írja le. A 3-helyes *olvas Lencsi könyvet úgyban* mondatváz esetében az eddig tárgyalt  $ie f$  szerkezet kiegészül egy  $-bAn$  raggal megjelölt helyel ( $g$ ), és a  $g$  helyet kitöltő *úgy* szóval ( $y$ ) a következő szerkezetté:  $ie f g$ . Ennek hálójá a 4. ábra mintájára képzelhető el, absztrakt szinten azonos a 3-dimenziós 3-Post-háló struktúrával. Az igei szerkezetek (mondatvázak) 3-Post-hálóra a következőkben az **isz-háló** (ti. igeiszerkezet-háló) rövidítést használjuk.

Vegyük észre, hogy egy olyan struktúrát kaptunk, ahol egy ontológiától vagy a WordNettől (Miháltz et al. 2008) vagy akár egy szóasszociációs hálózattól (Kovács et al. 2012) eltérően nem szavak, szócsoporthok, fogalmak vannak a csomópontokban, hanem szószerkezetek: különféle kitöltöttségű igei szerkezetek. Nem szavak vagy fogalmak relációit írják le ezek a hálók, hanem szerkezetek egymáshoz való viszonyait, abban az értelemben, hogy hogyan „hozható létre” az egyikből a másik. Az említett hálózatoktól eltérően ezekben a hálóknak a csomópontok szigorúan egymásra épülnek, közvetlenül feltételezik egymást.

Az éleken megjelennek olyan entitások – például az esetragok által képviselt helyek –, amelyek a felszínen nem képeznek közvetlen, önálló egységet, önálló szót. Az éleken lévő információegységekre – ahogy a 22. oldalon hívtuk őket – igaz az, hogy esetleg több ilyen egység is származhat egy felszíni szóból. Ha a kitöltők esetében egy pillanatra eltekintünk a szóösszetettől és az esetragokon kívüli toldalékoktól, és az *igazságérzet* kitöltőt egy (információ)egység lévén egy „morfémaként” kezeljük, azt mondhatjuk, hogy az éleken morfémák jelennek meg. A csomópontokban igei szerkezetek, az éleken pedig az igei szerkezet különböző szavaiból való morfémák vannak. Nem szóháló tehát ez, hanem a csomópontok tekintetében szerkezet-háló, az élek tekintetében pedig morfémaháló. A háló részbenrendezését és egyúttal szintezését az információegységek száma (*l*) adja meg: egy szerkezetből akkor mutat él egy másikba, ha az utóbbi egy információegységgel többet tartalmaz az előbbihez képest.

Felfogásunk annyiban egyetért az igék és vonzataik klasszikus valenciaalapú elméletével (Tesnière 2015 és Vincze 2011: 6. fejezet), hogy ige melletti helyekről, pozíciókról beszélünk, illetve hogy az alanyt ugyanolyan vonzatnak tekintjük, mint a többi vonzatot (Tesnière 2015: 51. fejezet). Abban viszont szembemegyünk a klasszikus felfogással, hogy „szabadon”, „igény szerint” teszünk hozzá plusz helyeket az igehez, nem vesszük eleve adottnak az ige meghatározott számú valenciáját. Másszóval: nem foglalkozunk azzal, hogy elvben milyen vonzatai vagy bővítményei lehetnek egy ige-nek, hanem azt nézzük korpuszban, korpuszvezérelt módon (Tognini-Bonelli 2001), hogy mi van mellette ténylegesen. Ezzel lényegében azt az álláspontot képviseljük, hogy eredendően érdemes az összes bővítményt egységesen kezelni (Čech et al. 2010).

## 5.2 A valódi igei szerkezet definíciója

A 6. ábrán látjuk, hogy egy mondatváz összes alszerkezete ténylegesen igei konstruktum: egy ígét és az ige bizonyos bővítményeit tartalmazza. Ezek olyan nyelvi egységek, melyek egy szótárba való bekerülésre esélyesek lehetnek. Azt fogjuk látni, hogy minden mondatváz alszerkezetei közül lesz egy, amit érdemes kiválasztani, kitüntetetten kezelni ebből a szempontból.

A korábban (15. oldal) bevezetett valódi szerkezet fogalomnak adunk most konkrét tartalmat az igei szerkezetek vonatkozásában. A **valódi igei szerkezet (v-isz)** a következő tulajdonságokkal rendelkezik: a vonzatokat megtestesítő esetragokat tartalmazza, a szabad határozókat megtestesítőket nem, ezenkívül az idiomatikus kitöltőket tartalmazza, az esetlegeseket (kompozicionálisakat) nem. A valódi igei szerkezet tehát *teljes*, azaz minden szükséges elemet tartalmaz, és *tiszta*, azaz semmilyen szükségtelen elemet nem tartalmaz. Az *olvas vmit v-isz*, az *olvas vmit vmiben* nem; a *vesz részt vmiben v-isz*, a *vesz részt* nem.

(Megjegyzés: esetraggal képviselt bővítményre jelen tanulmányban egyaránt használom a bevett *vmiben* típusú jelölést és a *-bAn* típusú jelölést is (az alany jele az utóbbi formában  $\emptyset$  lesz). Az utóbbi jelölés annyiban jobb, hogy nem utal az adott bővítmény élő/élettelen tulajdonságára, amivel a tanulmányban nem is foglalkozunk.)

A definíció belül a felhasználói célnak megfelelően rugalmasan eldönthető, hogy mit tekintünk vonzatnak, illetve idiomatikusnak. Mindenki alkalmazhatja a neki kedves vonzatsági és idiomatikusági meghatározást. Hasznosnak látszik – és jelen tanulmányban ezt fogjuk választani –, hogy érdemes ezeket tágran értelmezni. Az idiomatikuságnak bármilyen fajtáját, vagy akár „gyanúját” elegendőnek fogjuk tekinteni. Egy kitöltő a valódi igei szerkezet része lesz, ha az igevel nem kompozicionális vagy sajátos jelentéssel bíró vagy intézményesült (Sag et al. 2002) kifejezést alkot, vagy akár csak speciális módon fordítható egy másik nyelvre.

Arról van szó, hogy az ige melletti azon elemek együttesét keressük, melyek az igével együtt jelentésükben egységet (*unit of meaning* (Teubert 2005)) alkotnak: az igei szerkezet adott jelentésének megőrzéséhez minden elemre szükség van, azaz nem hagyható el elem a jelentés megváltozása nélkül; az igei szerkezet jelentéséhez szükséges minden elem megvan, azaz nem hiányzik elem; a meglévő elemek mással nem (vagy csak korlátozottan) helyettesíthetők; és csak megkötés nélkül (vagy legalábbis nagyméretű szóosztályból vett kitöltővel) kitölthető szabad helyek vannak benne.

Tekintve, hogy a szótár a nyelv jelentéssel bíró elemeinek gyűjteménye, ezek az önálló (nem kompozicionális) jelentéssel bíró igei szerkezetek pont a lexikográfiailag hasznos igei szerkezetek. Ezek azok, amiket egy szótárban szeretnénk látni. Már Kalmár László megállapítja egy elejtett megjegyzésében, hogy a szótárba az idiomatikus egységeket érdemes felvenni: „... a gépi fordítás céljára úgyis speciális szótárt (jobban mondva morféma- és idiomatárt) kell összeállítani.” (Kalmár 1964: 71) A szerzők nagyon egyetértenek azzal a felfogással, hogy a szótárba lényegében *csakis* az idiómák kelljenek, és a duplakocka modell segítségével a v-isz-ek számba vétele révén éppen ennek a célnak kívánnak megfelelni.

### 5.3 Két gondolat

Felvetődhet két gondolat, melyek látszólag alkalmas módon és sokkal egyszerűbben fogják meg az igei szerkezetek reprezentálásának problémáját. Mutassuk meg, hogy ezek a megközelítések nem lesznek nekünk elegendők, megfelelőek.

Az egyik gondolat az, hogy miért nem ábrázoljuk az igei szerkezeteket a hálós struktúra helyett  $R(\overset{\text{ige}}{\cdot}, \overset{\text{alany}}{\cdot}, \overset{\text{tárgy}}{\cdot})$  típusú relációkkal, ahol egyszerűen az előforduló konkrét igék, alanyok és tárgyak vannak egymással relációban.

Ez nem egy jó modell a szabad és kitöltött helyeket tetszőleges kombinációban tartalmazó isz-ek ábrázolására, hiszen ez pusztán a mondatvázak reprezentációja lenne. Ez a modell két nagyon fontos szempontot nem tud ábrázolni: (1) hogy egy isz-ben egyáltalán hány hely van, hány hely érdekes, és (2) hogy melyik helyen tekintendő az adott kitöltő az isz inherens részének, és melyiken nem. A hálós modell éppen ezeket képes kezelni: nem csupán egyes helyekre adja meg a jellegzetes kitöltőket, hanem azt is nézi, hogy egyáltalán melyik hely fontos/számít, illetve hogy hol kell, hol érdekes pusztán a hely; nem csupán azt mutatja meg, hogy mik a jellemző tárgyak, hanem hogy egyáltalán szokott-e tárgya lenni az adott igeinek. Ennek az alapja az, hogy a „viszonyokat” (helyeket) és a „dolgokat” (kitöltőket) elkülönítjük egymástól. A most említett modellt legfeljebb a duplakocka modell egy triviális kiindulópontjának tekinthetjük.

A másik gondolat, hogy ha már háló, akkor ugye az ige és a tklen szerkezetek (ige + helyek) által megadott „alsó egységkocka” (3. ábra, 17. oldal) lesz az „alapstruktúra” és a duplakocka többi része (a kitöltők) pedig a „realizáció”.

Nem! Ez pontosan az a realizáció-fogalom, amivel a jelen tanulmány élesen szemben áll, ami ellen küzd. Ezzel a felfogással éppen a lényegét, a *vesz részt vmiben* típusú vegyes szerkezeteket dobnánk ki. Éppen azt hangsúlyozzuk, hogy számos v-isz-nek elengedhetetlen része egy-egy kitöltő. Az a szándékunk, hogy semmiképp ne tekintsük a *vesz részt vmiben*-t a *vesz vmit vmiben* „realizációjának”. Ugyanis az utóbbi szerkezet nem létezik (értsd: nincs ilyen v-isz), az előbbi viszont egy teljes értékű, önálló v-isz. Másképp: az isz-hálók minden csomópontja teljes jogú isz, semmiképp nem oszthatjuk őket első- és másodosztályúakra, mindnek egyformán meg kell adnunk az elvi esélyt arra, hogy ő legyen az áhított v-isz. Abban az érte-

lemben persze beszélhetünk megvalósulásról, hogy egy adott mv-ban a megfelelő v-isz szabad helyei bizonyos kitöltőkkel ki vannak töltve, de ez egészen más, mint a most felvetett gondolat.

Vizsgáljuk most meg azokat a v-isz-eket, melyekben van kitöltő. Ige-névszó kollokációk lévén nevezhetjük őket **komplex igének** is. Ezek között van olyan, aminek van (az alanyon kívül) vonzata (*vesz részt vmiben, hűz hasznat vmiből* stb.), és van olyan is, aminek nincs (*kap észbe, jön létre, ér véget* stb.). Erős érvek szólnak amellett, hogy a komplex igéket önálló igéknek tekinthetjük: (1) a komplex igék az alapigétől független, önálló jelentéssel bírnak; (2) a komplex igéknek az alapigétől eltérő új vonzatkerete van (a *kap észbe* esetén eltűnik a tárgy, míg a *vesz részt vmiben* esetén egy új vonzat jelenik meg); (3) azonos alapigét tartalmazó komplex igék sok esetben eltérő alapigével fordíthatók más nyelvekre. Az egyes v-isz-eket tehát emiatt is önálló, egymás mellett létező, egymástól független entitásoknak tekintjük, akkor is, ha adott esetben két létező v-isz illeszkedik egymásra, mint ahogy ezt *vesz vmit* és a *vesz részt vmiben* esetében látjuk.

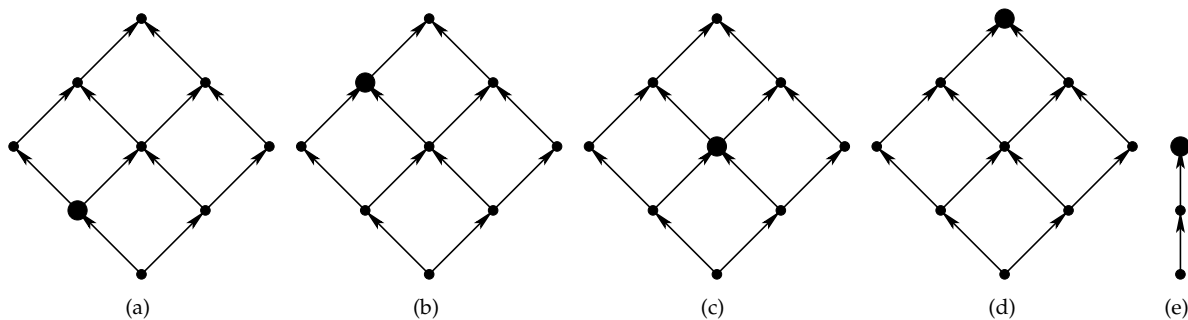
#### 5.4 Példák

Vegyünk néhány példát. Az említett *olvas Lencsi könyvet ágyban* mondatváz alszerkezetei közül a valódi igei szerkezet az *olvas vki vmit*. Az ebben a „mondatban” lévő valódi igei szerkezet egyszerűen a tárgyas *olvas* ige: a teljes szerkezet az ige mellett alanyt és tárgyat tartalmaz. A *-bAn* ragos helyhatározó nem része a valódi igei szerkezetnek, és nem részei a speciális jelentést nem hordozó, kompozicionális *Lencsi* és *könyv* kitöltők sem. A mondatváz és a v-isz egymáshoz való viszonyát a 7. ábra (a) részén látjuk. A szótárakban megszokott (Atkins és Rundell 2008) *tárgyas ige* megjelölés a modellünkben tehát így ragadható meg: ige + egy konkrét (tárgyi) hely. Ez egy példa arra, hogy az ismert, megszokott fogalmak a modellbe beleilleszkednek, a modell alkalmas azok megragadására. Alapvető szótári hagyomány ilyen „többinformációegységes” entításokat szerepeltetni a szótárban. A tárgyas ige éppen arra példa a hagyományból, aminek az általánosítását a hálós modell valósítja meg.

A *vesz Lencsi részt felolvasáson* mondatváz valódi igei szerkezete a *vesz vki részt vmin*. Az ige mellett a szerkezet része a tárgy a konkrét tárgyi kitöltővel együtt, az alany és egy *-n* ragos vonzat. A ’participate in’ jelentéshez nyilván szükséges a tárgyi helyet kitöltő *rész* szó jelenléte, és szükséges az *-n* ragos vonzat, ami megadja, hogy a részvétel mire irányul. Nem szükséges viszont a speciális jelentést nem hordozó, más szavakkal helyettesíthető *Lencsi* és *felolvasás* kitöltő (7. ábra (b)).

Említettük, hogy a valódi igei szerkezetekből a jelentés megőrzése mellett nem hagyható el elem. Természetesen a *Lencsi* vagy a *felolvasás* elhagyásával is „változik” a jelentés (szorosán véve *mást* fog jelenteni az adott isz): de azt a különbséget kell itt látni, ami a *Lencsi*, illetve a *rész* elhagyása között van, másként megközelítve ami a *Lencsi Mici*-re vagy *Csöpi*-re való cserélése, illetve a *rész elégtétel*-re vagy *autó*-ra való cserélése között van. Az utóbbi az amit, egy v-isz esetén a v-isz sérülése nélkül nem tehetünk meg. A v-isz részét képező elem elhagyása/helyettesítése az egész jelentést befolyásolja, a v-isz részét nem képező elem elhagyása/helyettesítése csak a saját hozzáadott jelentését érinti.

A v-isz-ek változatos kombinációkban tartalmazhatnak szabad és kitöltött helyeket: *kerül sor vmire* – egy szabad és egy kitöltött hely (7. ábra (b)), a *ad vki vkinek vmit* – három szabad hely (7. ábra (c)), a *fest vki ördögöt falra* – egy szabad és két kitöltött hely (7. ábra (d)), a *Kisütött a nap*. mondat esetében pedig nincs szabad hely, a v-isz azonos a mondatvázal (7. ábra (e)).



7. ábra. Valódi igei szerkezetek elhelyezkedése mondatvázak hálóján. A ábrán legfeljebb két dimenziót tüntetünk fel, a kitöltött alanyt tartalmazó v-isz-ek – a (b) második példája és az (e) – kivételével az alanyi dimenziót elhagyjuk! (a) Az *olvas vmit* v-isz az *olvas könyvet ágyban* mondatváz hálóján. (b) A *vesz részt vmin* v-isz a *vesz részt felolvasáson* mondatváz hálóján. Ugyanez az ábra jeleníti meg a *kerül sor vmire* v-isz-t, csak persze tárgy és *-n* helyett alany és *-ra* a két hely. (c) Az *ad vkinek vmit* v-isz az *ad Jáninak könyvet* mondatváz hálóján. (d) A *fest ördögöt falra* v-isz a *Ne fessd az ördögöt a falra!* mondat mondatvázának hálóján. (e) A *kisüt nap* v-isz a *Kisütött a nap.* mondat mondatvázának hálóján.

Látjuk, hogy a hálóban a v-isz a mondatvázhoz képest bárhol elhelyezkedhet, azaz adott esetben bármelyik alszerkezet lehet v-isz, akár a maximális szerkezet vagy a gyökér is. A példák alapján az azonban kézenfekvőnek látszik, hogy egy mondatváz alszerkezetei között *egy darab* valódi igei szerkezet van, azaz mindig egyértelműen eldönthető, hogy egy adott mondatváz esetében mi a v-isz. Ez triviálisnak tűnik, mindenesetre a továbbiakban feltesszük, hogy így van. Ha esetleg konkrét esetben több v-isz-jelölt is felmerülne, akkor közülük azt tekintjük v-isz-nek, ami a leginkább megfelel a definíciónak.

Amint látni fogjuk, az adott (tag)mondatokhoz tartozó v-isz-ek beazonosítása, megtalálása lehet az egyik célunk. Ezt a dolgot későbbi részében (7. rész, 34. oldal) érintjük.

## 6 A korpuszháló felépítése

Ebben a fejezetben szintet lépünk. Eddig a modell alapjait fektettük le részleteiben kidolgozva, viszonylag egyszerű fogalmakat vezettünk be, egyszerű megállapításokat tettünk. Eddig tagmondatokat reprezentáltunk, a nagy lépés abban rejlik, hogy most teljes korpuszt fogunk a modell segítségével ábrázolni. Tegyük fel, hogy rendelkezésünkre áll egy nagy méretű korpusz valamiféle (automatikus) szintaktikai elemzéssel, amiből megállapíthatók az ige és a bővítmények szükséges jellemzői. Az ötlet az, hogy építsük fel valamilyen módon az egész korpusz hálóját, azaz egy olyan struktúrát, ami a korpusz összes igei szerkezetét magában foglalja, reprezentálja. A 2-4. fejezetekben követett absztrakt tárgyalásmód itt háttérbe szorul, mindig igei szerkezetekről fogunk beszélni, de természetesen az absztrakt modell is mindig érvényes és hozzágondolható.

A teljes korpusz alapján készített összetett hálórendszer **korpuszháló**nak fogjuk nevezni. A korpusz összes mondatvázának hálójából, e háló integrálásával fogjuk felépíteni. Eddig alapelemekből építettünk isz-hálókat, most az isz-háló lesznek ennek az újabb építkezésnek az építőkövei.

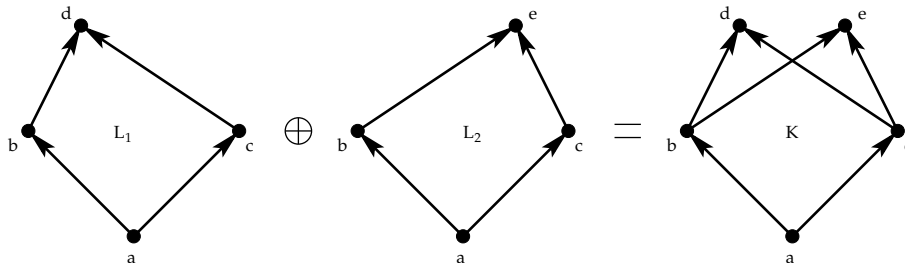


### 6.1 Hálók beágyazása és kombinálása

Szükségünk lesz egy olyan műveletre, mely két isz-hálót egy struktúrává kombinál össze. Ehhez bevezetünk néhány fogalmat. Egy  $L$  háló **részhálója**  $K$ -nak (jele:  $L \subseteq K$ ), ha  $L$  csomópontjai és élei  $K$ -ban is megvannak, és  $L$  maga is háló. Egy  $L$  háló **beágyazható** egy (nagyobb)  $K$  hálóba, ha létezik izomorfizmus  $L$ -ről  $K$  egy részhálójára (Partee et al. 1990: 287). Ez szemléletesen annyit jelent, hogy a  $K$  rajza tartalmazza  $L$  rajzát. A hálónál (ld. 16. oldal) gyengébb fogalom a (véges) **félháló**: csak azt követeljük meg, hogy *vagy* egyértelmű minimuma *vagy* egyértelmű maximuma legyen. Az előbbi a  $\wedge$ -félháló (meet-félháló), az utóbbi a  $\vee$ -félháló (join-félháló). Szemléletesen: ha vízszintesen középen kétfelé vágjuk a 2. ábrán (16. oldal) látható hálót, a  $bac$  egy  $\wedge$ -félháló lesz a  $bdc$  pedig egy  $\vee$ -félháló. Félhálókra a részfélháló és a beágyazás a fentiekhez hasonlóan értelemszerűen definiálható.

A matematikában sokszor címkézés nélkül tekintik a struktúrákat. Az isz-hálók esetében viszont nyilván kiemelten fontos a címkézés, azaz hogy a csomópontokban konkrétan milyen szerkezetek, és az éleken milyen alapelemek és műveletek vannak. A különböző módon címkézett de azonos felépítésű isz-hálók élesen különböznek. E különbségtétel érdekében minden isz-háléhoz egy címkézőfüggvényt fogunk hozzáérteni (ahogy impliciten eddig is tettük), amely a háló csomópontjaihoz és éleihez hozzárendeli a megfelelő entitásokat. Ennek fényében a fent emlegetett izomorfizmus valójában azonosság.

Most meghatározzuk a **hálók kombinálása** (jele:  $\oplus$ ) műveletet, mely az azonos gyökérrel bíró hálók kombinálására szolgál. Legyen  $L_1 \oplus L_2 = K$  úgy, hogy  $K$  címkézéshelyes, minimális  $\wedge$ -félháló, amibe mindkét háló beágyazható. Másképp: legyen igaz, hogy  $L_1 \subseteq K$  (címkézéshelyesen) és  $L_2 \subseteq K$  (címkézéshelyesen) és  $K$  a lehető legkevesebb csomóponttal és éllel bír, és  $L_1 \cup L_2$  címkézett pontjai mind csak egyszer szerepelnek  $K$ -ban (8. ábra).



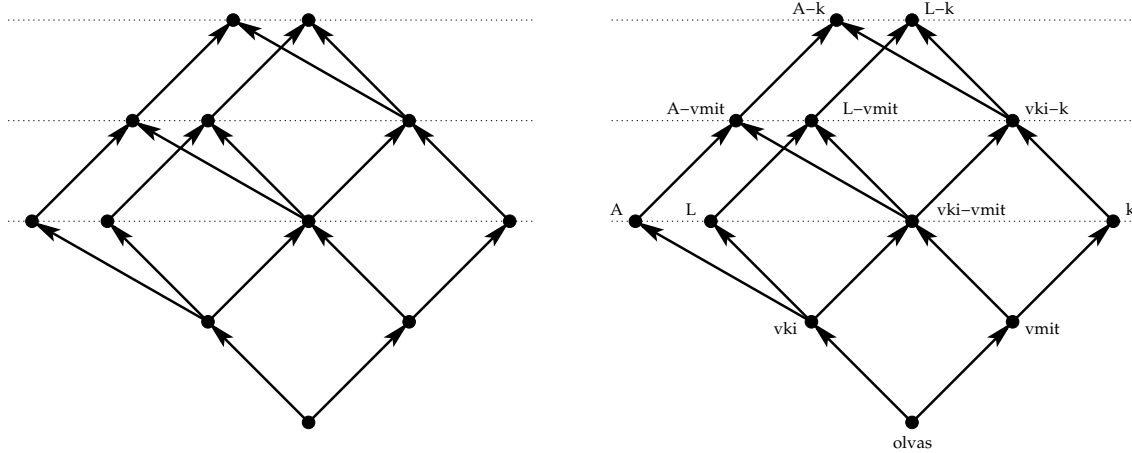
8. ábra. A  $\oplus$  (hálók kombinálása) művelet illusztrációja. A beágyazás eredményeképpen  $L_1$  és  $L_2$  rajza részét képezi  $K$  rajzának.  $K$   $\wedge$ -félháló: egyértelmű minimuma van, egyértelmű maximuma nincs.  $K$  minimális és címkézéshelyes. Ha címkézés nélkül tekintenénk a hálókat, akkor  $L_1 \equiv L_2$  lenne, sőt a kombinálás eredménye is ugyanez a háló lenne. Szemléletesen ez történik kombináláskor: ahol a két kiinduló hálónak metszete van, ott egybeesik, ahol nincs, ott elágazik az eredményül kapott  $\wedge$ -félháló. (Ezen az ábrán a könnyebb érthetőség kedvéért nem isz-háló – nem duplakocka –, hanem egyszerűbb hálók szerepelnek.)

Csak azonos gyökérrel bíró hálókat (azonos igéhez tartozó mondatvázakat) fogunk kombinálni, ezért a  $\oplus$  művelet eredménye mindig létezni fog, valóban mindig félháló – egész pontosan  $\wedge$ -félháló – lesz. A megkövetelt közös egyértelmű minimális elem a gyökér, azaz éppen az ige lesz. Amikor felépítjük a teljes korpuszhálót, akkor valójában külön-külön igénként építünk egy nagy méretű  $\wedge$ -félhálót, ezért nevezhetjük a korpuszhálót hálórendszernek. Sok esetben a korpuszhálónak csak egy adott igt tartalmazó részével fogunk foglalkozni, az egyszerűség kedvéért ezt is korpuszhálónak fogjuk nevezni.

## 6.2 Példák

Nézzünk néhány példát a fent definiált  $\oplus$  művelettel előállított (egyszerűbb) kombinált  $\wedge$ -félhálóra, azaz olyan struktúrákra, ahol több duplakocka ábrázolódik egyben: többféle kitöltő, többféle hely fordul elő. Ezeket a kombinált  $\wedge$ -félháló elnevezés mellett szemléletesen **bo-kornak** is fogjuk hívni.

A legegyszerűbb eset, mikor két olyan mondatvázat kombinálunk, ami egyetlen kitöltőben tér el egymástól (9. ábra). Ezeken az ábrákon a kitöltők kezdőbetűit tüntetem csak fel, kötőjellel összekapcsolva, a helyeknek az első szinten ( $l = 1$ ) látható sorrendjében. Az  $A$ -k tehát adott esetben azt jelenti, hogy *Anya könyvet*.



9. ábra. *Lencsi könyvet olvas.*  $\oplus$  *Anya könyvet olvas.*  $\equiv i e f \oplus i e f$ . A várt „közös alsó rész, eltérő felső rész”  
 $w x \quad w y$   
 struktúrát látjuk. Így képzelhetjük el: egy ponton/síkon elvágjuk a duplakockát, egy elágazást csinálunk, és onnantól még egyszer ( $n$ -szer ha  $n$  ilyen mondatvázat kombinálunk) felépítjük az elvágott részt, ami ezáltal az eredetihez képes „felnyílik”. Ez hasonlóan működik  $h > 2$  esetén is.

A következő példában három mondatvázat kombinálunk, a második, illetve a harmadik más-más helyen bár, de szintén csak egy kitöltőben tér el az elsőtől (10. ábra).

A harmadik példa csak annyiban más, hogy megengedjük, hogy a két helyen egyszerre térjen el a két kitöltő. Ez négy mondatvázat jelent (11. ábra).

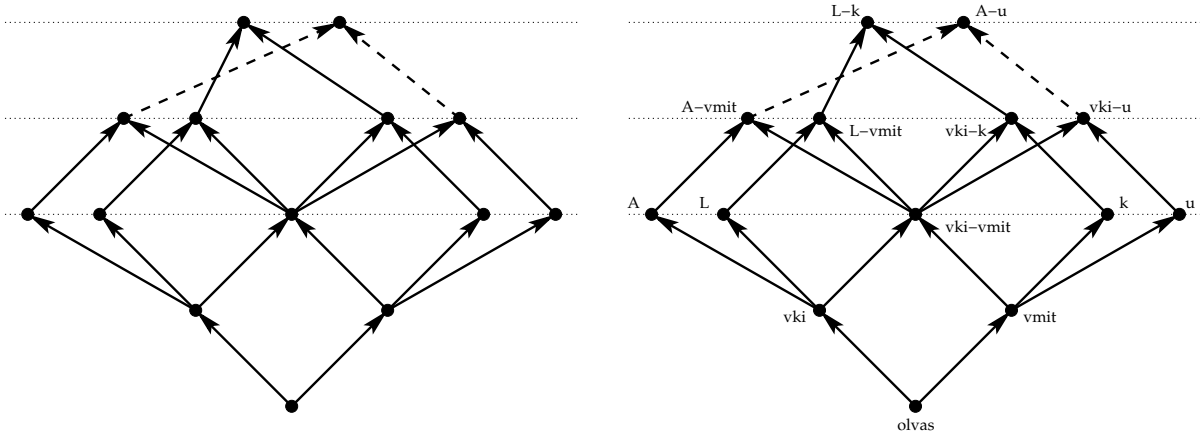
Az utolsó példában ismét csak két mondatváz szerepel. Itt két helyen egyszerre tér el a kitöltő (12. ábra).

A fentiekben csak azonos dimenziószámú ( $h = 2$ ) Post-hálókat kombináltunk. Hogyan néz ki ez különböző dimenziószámú hálók esetén? Mindig két lehetőség van: vagy egybeesik a két háló bizonyos része, vagy nem. Egy kisebb dimenziós (de beágyazható) duplakocka például rá fog simulni a nagyobb dimenziós háló egyik oldalára, ahogy azt a 17. oldalon láttuk. A nem illeszkedő részek elágazó, szélső esetben akár a gyökértől kezdve teljesen külön álló ágakat fognak létrehozni. Utóbbi akkor fordulhat elő, ha egy hely (eset) sose kombinálódik másokkal, mindig csak önmagában, a többi esettől elkülönülten szerepel az adott ige mellett.

## 6.3 Metaháló és koszorzat

Érdeemes a fentieket kategóriaelméleti keretben (Mac Lane 1998) is áttekinteni. Azt fogjuk látni, hogy a  $\oplus$  művelet a kategóriaelmélet egyik standard műveleteként értelmezhető.





12. ábra. *Lencsi könyvet olvas.*  $\oplus$  *Anyu újságot olvas.*  $\equiv i e f \oplus i e f$ . Ezt az ábrát a 11. ábrából a két szélső maximális szerkezet elhagyásával kapjuk. Az előző ábráktól eltérően itt két olyan maximális szerkezetet látunk, melyeknek csak két szinttel lejjebb van közös pontjuk, mégpedig a tklen szerkezet. Azaz a struktúra élesen eltér a 9. ábrán láthatótól. Így képzelhetjük el: itt egy  $L$  alakú „ollóval” vágjuk el a duplakockát, és a vágási vonaltól építjük fel a felső részt még egyszer (illetve annyiszor, ahány  $ef$ -be szánt párunk van). A struktúra emlékeztethet a Rubik-féle *Karikavarázs* játéokra.

Az isz-hálókon és az isz-hálókból kombinálással képzett  $\wedge$ -félhálókon (bokrokon) kívül képzeljünk el egy harmadik fajta hálószerű struktúrát. Ennek objektumai a kombinált címkézett  $\wedge$ -félháló lesznek, és akkor lesz közöttük  $A \rightarrow B$  nyíl, ha  $A$  beágyazható  $B$ -be. Ez a struktúra bizonyos értelemben „hálók hálója”, ezért **metaháló**nak nevezzük. Mivel  $(L \oplus L = L)$  révén az isz-háló is tekinthető kombinált  $\wedge$ -félhálónak, a metaháló objektumai között természetesen megtaláljuk az isz-hálókat is. A metaháló csúcsa, maximális eleme pedig éppen a korpuszháló lesz.

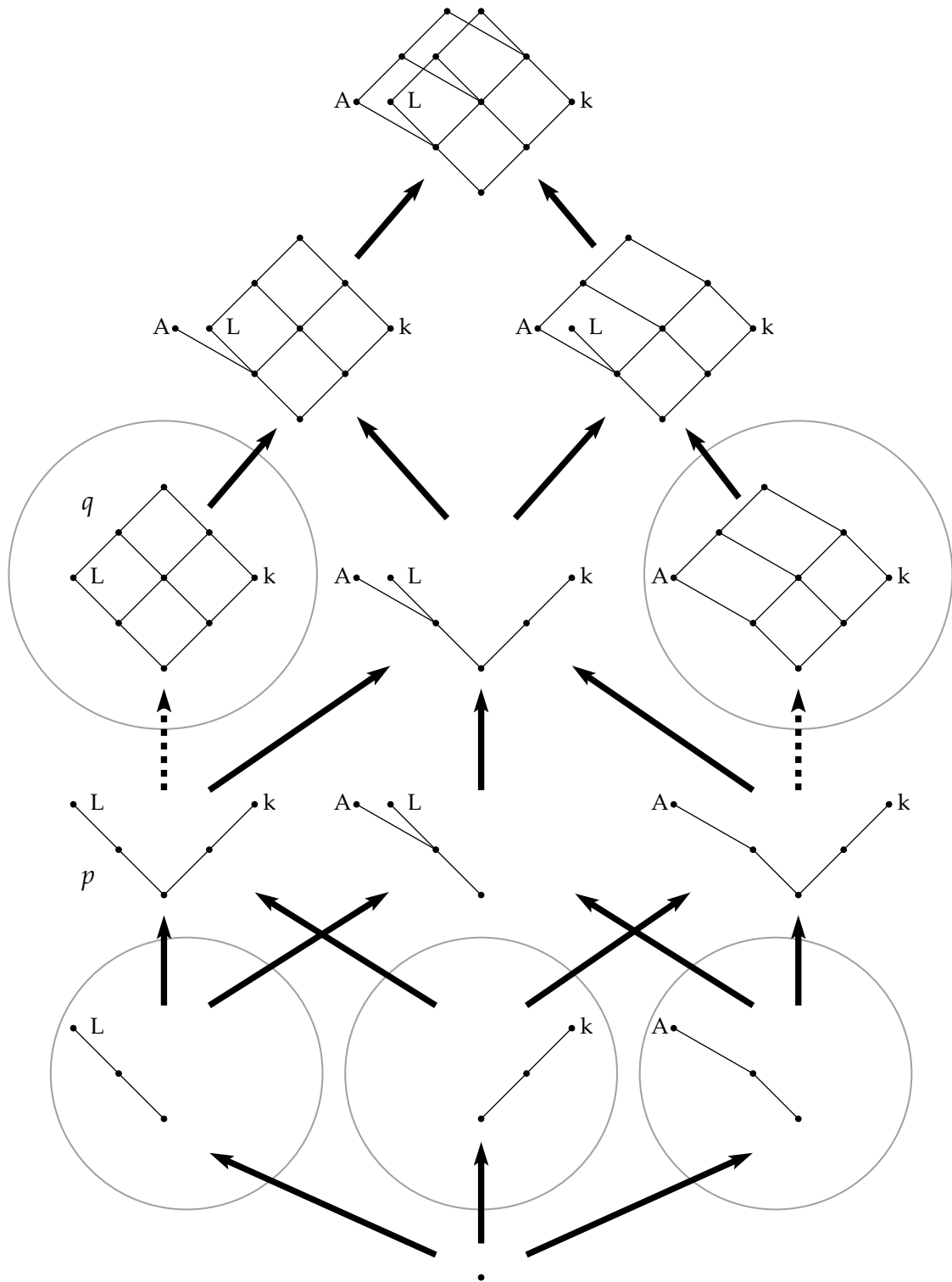
Definiáljuk a **koszorzat** (másnéven kategoriális összeg) fogalmát. Az  $N_1$  és  $N_2$  objektum koszorzata az az  $M$  objektum, amely mindkettőből nyilakon keresztül elérhető (kategoriáleméleti megfogalmazásban:  $\exists N_1 \rightarrow M$  és  $\exists N_2 \rightarrow M$  morfizmus), és az ilyenek közül minimális, azaz az ilyenek összességén belül nincs feléje mutató nyíl (Mac Lane 1998: 62-63). (Jó szemléletes definíciót ad még Steve Easterbrook előadásának 25. diája

<http://www.cs.toronto.edu/~sme/presentations/cat101.pdf>) A koszorzat a háló „nyilain fölfelé haladó” maximumképzés ( $\vee$ , join) művelet általánosítása.

A fentiekből adódik az állítás, hogy az imént bevezetett hálók kombinálása művelet éppen a metaháló struktúrán értelmezett koszorzattal azonos.

A háromféle egyre bonyolódó struktúra – az isz-háló, a bokrok és a metaháló – nem triviális viszonyát egy alkalmas hasonlattal lehet megvilágítani. A hasonlat forrástartománya a természetes számok halmaza lesz, a szokásos szorzás művelettel, prímszámokkal, összetett számokkal és oszthatósággal. A csak gyökekből álló isz-háló az egységelem. Az isz-háló az építőkövek – hasonlóan a prímekekhez. A bokrok ezekből az építőkövekből jönnek létre – hasonlóan az összetett számokhoz. A létrehozó  $\oplus$  művelet a szokásos szorzásnak felel meg, a bennfoglaló algebrai struktúra – a metaháló – a természetes számok halmazának, a *beágyazható* reláció pedig az *osztója* relációnak.

A 9. ábra metahálója, pontosabban a 9. ábrán megtalálható 0, 1 és 2 dimenziós duplakockák metahálója a 13. ábrán látható.



13. ábra. Egy metaháló-részlet: a 9. ábrán (és jelen ábra csúcsán) látható struktúrát létrehozó  $\wedge$ -félhálók metahálója. A bekarikázott entitások a prímek (isz-hálók), a nem bekarikázottak az összetettek (bokrok), a vastag nyilak az *osztója* (*beágyazható*) relációt jelölik. Tetszőleges két  $\wedge$ -félháló kombinációja (koszorzata) a nyilak mentén fölfelé haladva az legközelebbi közös pontban található. Egy  $\wedge$ -félháló rajzán belül a vonalak pontos irányának persze nincs jelentősége, de a könnyebb áttekinthetőség kedvéért mindig ugyanúgy vannak megjelenítve, illetve ugyane célból a címkek még tovább rövidítve szerepelnek.

A sok hasonlóság mellett a metahálónak van egy érdekes tulajdonsága, ami eltér a természetes számok struktúrájától. Vizsgáljuk meg a 13. ábra bal oldalán a  $p$  és a  $q$  bokrot, és az őket összekötő szaggatott nyilat.  $q$ -nak van önmagától és az egységtől különböző osztója (például  $p$ ), ezért  $q$  nem lehet prím. Látjuk, hogy az ábra nem bekarikázott  $\wedge$ -félhálói összetettek: mind létrehozhatók saját maguktól különböző osztóikból a  $\oplus$  művelettel, alkalmasan választott vastag nyilak mentén. Vegyük észre, hogy  $q$  viszont nem hozható létre (nála kisebb) osztóiból a  $\oplus$  művelettel, azaz ebben az értelemben nem összetett. Se nem prím, se nem összetett, mi lehet akkor?

Nézzük meg jobban a kombinálás műveletet. A metahálóban nincs „rendes” hatványozás, a kombinálás idempotens:  $\forall s : s \oplus s = s$ , azaz egy  $\wedge$ -félhálót saját magával összekombinálva ismét csak saját magát kapjuk. Vigyük át ezt a tulajdonságot a hasonlat másik oldalára, a természetes számok körére. Minden marad a régiben, szorzás, prímekek, oszthatóság, de vezessünk be egy új szabályt, miszerint:  $n \cdot n = n$ . Ekkor (bizonyítható, hogy) a szorzás és a legkisebb közös többszörös (LKKT) művelete ugyanazt az eredményt adja, a két művelet egybeesik. Azaz a kombinálásnak nemcsak a szorzás, hanem az LKKT is egy analógiája a hasonlat másik oldalán.

Szorzás helyett LKKT-t véve a természetes számok között van egy csoport, ami rendelkezik a fentebbi tulajdonságokkal, ti. hogy van önmagától és az egységtől különböző osztója, ugyanakkor nem hozható létre nála kisebb osztói LKKT-jaként. Ezek a prímhatványok:  $\{x^y : x \in \mathbb{P} \wedge y \in \mathbb{N} \wedge y \geq 2\}$ . A metahálóban azonban az idempotencia miatt a prímhatványok nem jelennek meg külön, hanem egybeesnek a prímekekkel. Így  $q$  a szó szoros értelmében nem nevezhető prímhatványnak sem, pedig a tulajdonságai efelé mutatnak.

A fentiek alapján a  $q$ -t (és hozzá hasonlóan az összes magasabb dimenziós isz-hálót) a prímhatványokra emlékeztető viselkedésük miatt külön terminussal **magasabb dimenziós prímekeknek**, illetve az isz-háló dimenziójának megfelelően  **$n$ -dimenziós prímekeknek** nevezhetjük. Ezeknek csak alacsonyabb dimenziós osztóik vannak, a saját dimenziójukon belül tehát valóban prímként viselkednek. Fent mondtuk, hogy az isz-hálók prímekek, ez továbbra is így van, csak a metaháló prímfogalmát bontottuk ki egy kicsit.

Megjegyezzük, hogy a  $\oplus$  művelet kommutatív, asszociatív és nem invertálható. A metaháló csomópontjai tehát kommutatív félcsoportot alkotnak a  $\oplus$  művelettel, éppen mint a természetes számok halmaza a szorzás művelettel. Talán nem véletlen, hogy Mazur (2008: 10. rész) is a természetes számokban találja meg egy fontos absztrakció kiindulópontját.

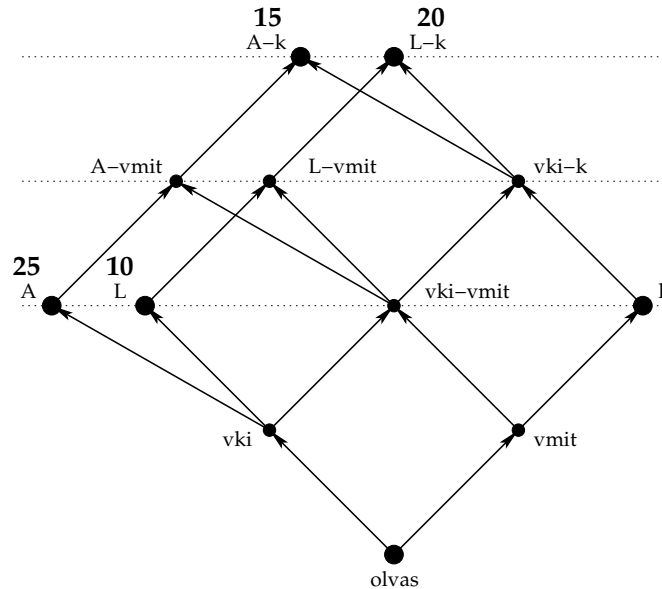
#### 6.4 A korpuszgyakoriságok szerepe

Egy korpuszban sok olyan mondatváz van, ami többször is előfordul. Szigorúan a fentiek alapján építve a korpuszhálót ezek csak egyszer jelennének meg. Annak érdekében, hogy a jelentőségüket megragadjuk, hogy az egyszer és az ezerszer előforduló mondatváz ne azonos módon szerepeljen a korpuszhálóban, azt a megoldást választjuk, hogy a korpuszháló megfelelő pontján feljegyezzük az adott mondatváz gyakoriságát. Másképp: minden mondatvázhoz rendelünk egy számlálót a mondatváz maximális szerkezetének megfelelő ponton, ami a mondatváz előfordulásait számolja. Ahogy a korpuszhálót építve haladunk végig a korpuszon, és olyan mondatvázhoz érünk, ami már szerepelt, annyit teszünk csak, hogy eggyel megnöveljük a hozzá tartozó számlálót. Így a kész korpuszháló a mondatvázakhoz rendelt gyakorisági értékekkel is fel lesz szerelve.

A korpuszban nyilván minden helyen van valamilyen kitöltő, így korpuszgyakorisági értékeink (jele a továbbiakban:  $f$ ) csak tk csomópontoknál lesznek. A 3. ábrán (17. oldal) lévő

hálóban a nagy korongok jelölik azokat a csomópontokat, melyekhez gyakorisági számláló van rendelve. Ezen a korábbi ábrán látható háló 4 darab (azaz  $2^h$ , vö: 2. táblázat) tk csomópontot tartalmaz: van benne egy 2-dimenziós, két 1-dimenziós és egy 0-dimenziós mondatváz.

Illusztrációképpen a 14. ábrán látható az előző oldalacról már ismert bokor (fiktív) gyakorisági adatokkal ellátva.



14. ábra. Fiktív gyakorisági adatokkal ellátott bokor. Gyakorisági adat a nagy korongoknál szerepelhet. Plauzibilis, hogy egy korpuszban nem fordul elő az *olvas* (legalább az igeragból kikövetkeztethető) alany nélkül (a gyökér ill. a jobb oldali nagy korong), viszont tárgy nélkül előfordulhat (25 és 10). Természetesen szerepelhet a struktúra alsóbb szintjén kisebb szám: esetünkben 10 olyan mondatváz van, ahol *Lencsi* az alany és nincs tárgy, ezen kívül 20 olyan mondatváz, ahol a *Lencsi* alany mellett a *könyv* tárgy is szerepel. (Azok a tk szerkezetek, ahol nincs gyakorisági adat (értsd: nulla), azok csak a nagyobb dimenziós mondatvázak alszerkezeteként jelennek meg.)

## 7 A valódi igei szerkezetek meghatározása

### 7.1 Csomósodás és globális valódi igei szerkezetek

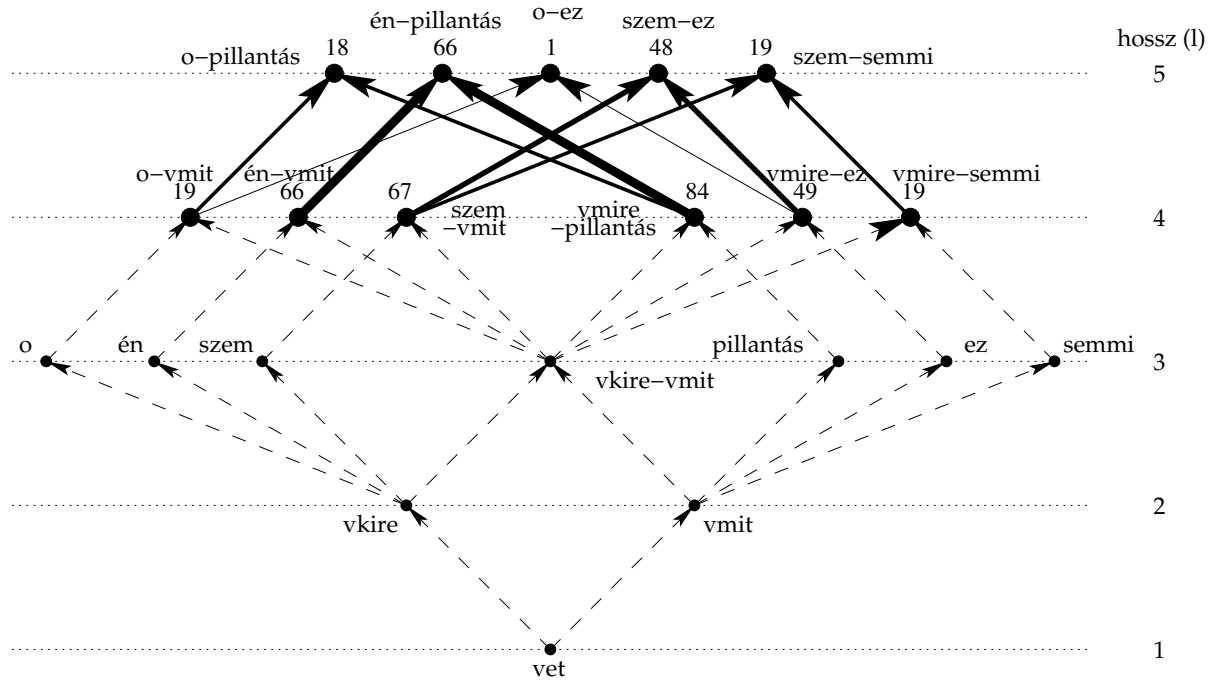
Az igei szerkezetek algebrai modellje ismertetésének – egyúttal a tanulmány lényegi részének – végére értünk. Mire lehet jó ez a modell, azon kívül, hogy új szemléletet ad? Milyen gyakorlati haszna, alkalmazása lehet? Alkalmas lehet például a valódi igei szerkezetek megragadására. Egy erre szolgáló módszer elvi vázlatával foglalkozunk a záró fejezetben.

Vegyünk egy valóshoz közeli összetettebb példát, mely egy gyakorisági adatokkal ellátott igazi korpuszháló egy kis részletének tekinthető, és vizsgáljuk meg a tulajdonságait.

Az alábbiakban a már említett *vet pillantást vmire* versus *vet vmit szemére* példát fogjuk használni. Fontos a két szerkezet viszonya: mindkettőben ugyanazt a két helyet látjuk, de az elsőben a tárgyi hely van kitöltve és a *-ra/-re* ragos kitöltetlen, a másodikban pedig fordítva (Sass 2011: 5,76-77).

A szabadon hozzáférhető Mazsola adatbázisból (Sass 2015) válogattuk ki a *vet* ige mellett tárgyi és *-ra/-re* ragos helyet tartalmazó mondatvázakat, hogy az illusztráció reális gyakorisági

értékeket tartalmazzon. A gyakorisági értékeket a példa kedvéért a többi helytől eltekintve összesítettük, és mindkét hely esetében plusz 2-2 kitöltőt választottunk ki a leggyakoribbak közül, így a következő szerkezeteink vannak: *vet + pillantást/ezt/semmit + szemére/rám/rá*. A szóban forgó mondatvázaknak a  $\oplus$  művelet segítségével egyesített hálóját a 15. ábrán látható.



15. ábra. Az „egész korpuszból” felépített korpuszháló illusztrációja, a teljes korpuszháló egy kis részlete. Az ábra a *vet + pillantást/ezt/semmit + szemére/rám/rá* összes kombinációjából adódó 9 db mondatváz közül azt az ötöt jeleníti meg, amely előfordul a korpuszban. Az ábra a 11. ábrával rokon. A csomópontokban a kitöltők szótóvét tüntettük fel, kötőjellel összekapcsolva, a helyeknek az  $l = 2$  szinten látható sorrendje szerint: az *én-pillantás* olvasata tehát *rám pillantást*. Formálisan a  $\oplus_{w \in \{\text{szem, én, ő}\}, x \in \{\text{pillantás, ez, semmi}\}} i e f_{w x}$  bokorról van szó, ahol  $i = \text{vet}$ ,  $e = -ra/-re$  és  $f = -t$ . Az ábrán a korpuszgyakoriságok természetesen a felső ( $l = 5$ ) szinten szerepelnek. Vegyük észre, hogy az  $l = 4$  szinten nem tk szerkezetek vannak, itt nem kellene megjelenniük gyakorisági értékeknek. Ezek az értékek az  $l = 5$  szintről *összesített* gyakoriságokat mutatják (erre még visszatérünk), a felső két szint közötti nyilak vastagsága az  $l = 5$  szintről származó érték nagyságát jeleníti meg. Persze a valóságban mindkét ábrázolt helyen számos további kitöltő fordul elő, az ábra mindkét oldalán hozzáképzeltük az ezeknek megfelelő hálórészleteket.

Nyilván az lenne a kívánatos, ha esetünkben a *vet pillantást vmire* és a *vet vmit szemére* lenne a v-isz. Legalábbis a 15. ábra első két mondatváza tekintetében az előbbi, a 4. és 5. mondatváza tekintetében pedig az utóbbi. (A 3. mondatváz esetében vélhetően egy rövidebb szerkezet lesz a v-isz, talán egyszerűen a *vet vmit*.) Azaz az lenne a jó, ha az azonos két helyet tartalmazó mondatvázak adataiból valamilyen (automatikus) eljárás alapján rá tudnánk jönni, hogy az egyik v-isz esetén az egyik hely van kitöltve, a másik esetén pedig a másik.

Ha jobban megvizsgáljuk, azt látjuk, hogy a megoldást jelentő  $l = 4$  szintű csomópontok bizonyos közös tulajdonságokkal bírnak. Úgy tűnik, hogy a v-isz-t képviselő csomópontok a korpuszhálónak valamiféle csomósodási pontjaiban vannak: olyan pontokban, melyek több, nagy gyakorisági értékű ponttal állnak kapcsolatban. És mindkét szempont fontosnak tűnik itt: hogy nagy gyakorisági értékű pontokkal, és hogy több ilyen ponttal. A példán azt látjuk, hogy az e két tulajdonságnak megfelelő csomópontok helyesen jelölik ki a v-isz-eket.



Fontos megjegyezni, hogy a fent felvillantott módon nem a **lokális v-isz**-eket (ld. 24. oldal), hanem a **globális v-isz**-eket ragadjuk meg. Azaz nem azt, hogy egyes konkrét mondatokban mi a v-isz, hanem hogy a teljes korpuszban mely szerkezetek szoktak leginkább v-isz-ek lenni. Másképp fogalmazva egy nyelv – jelen esetben a magyar – jellegzetes v-isz-eit próbáljuk megkeresni, felfedezni.

A globális v-isz lista, v-isz „tár” birtokában aztán lehet találni olyan módszert, ami a lokális v-isz-ek meghatározásában segít, ami a globális v-isz-eket visszavetíti az egyes mondatokra. Lényegében annyit kell tenni, hogy a mondat mondatvázára illeszkedő (ld. 15. oldal) globális v-isz-ek közül ki kell választanunk az igazit. Kézenfekvő kiinduló ötlet lehet, hogy a legspecifikusabbat (azaz a leghosszabbat) választjuk. Ez az esetek jelentős részében várhatóan elég jól meg fog felelni a lokális v-isz-nek. E kérdés további aspektusaival jelen tanulmányban nem foglalkozunk.

## 7.2 Örököltetés és az algoritmus kerete

Mostanra világos, hogy hogyan tudunk felépíteni egy teljes korpuszhálót a gyakoriságokkal együtt. Ahogy az előző részben már utaltunk rá, a feladat most az, hogy a korpuszháló segítségével meghatározzuk a valódi igei szerkezeteket. Az előbbieken körvonalazott módszerhez teszünk hozzá további megfontolásokat. Vizsgáljuk tovább a 15. ábrát.

Formailag minden szerkezet minden egyes nála (a háló részbenrendezése értelmében) kisebb (vö: 2. ábra a 16. oldalon) szerkezetet megvalósít. Más szóval a rá illeszkedő szerkezeteket. Az *olvas könyvet ágyban* (a továbbiakban:  $o_1$ ) szerkezetben rejlik egy *olvas vmit vmiben* ( $o_3$ ) szerkezet, hasonlóan ahogy a *vesz részt vitában* ( $v_1$ ) szerkezetben egy *vesz vmit vmiben* ( $v_3$ ).

Azt fogjuk mondani, hogy a gyakoriságokat egy szerkezetre **örököltetjük**, mikor az adott szerkezethez hozzárendelünk egy ún. **gyakorisági mérőszámot** ( $q$ ), melyet a nála (hálói értelemben) 1-gyel nagyobb szerkezetek gyakorisági mérőszámaiból képezünk valamilyen örököltetési módszer ( $I$ ) alkalmazásával.

$$q(x) = I(y_1, \dots, y_k) : x \leq y_i \wedge l(x) + 1 = l(y_i)$$

Az örököltetés tehát a szokásos nyilakkal szemben, „lefelé” történik, a kevés számú (vö: 2. táblázat, 20. oldal) „1-gyel kisebb” csomópont felé. A gyakorisági mérőszám kiinduló értékét a teljesen kitöltött szerkezeteknél adjuk meg, és az az egyszerű korpuszgyakorisággal lesz azonos (ld. a 15. ábra  $l = 5$  szintjét). És utána a korpuszháló tetejétől kezdjük el az örököltetést, szintenként haladva. A rövidebb tk szerkezeteknél, ahol korpuszgyakorisági érték is szerepel, ehhez adódik hozzá a kiszámolt  $q$ .

Ún. **összegzőes** örököltetés esetén az örököltetési módszer az, hogy a gyakorisági értékeket minden lehetséges „1-gyel kisebb” (vö: 3.2. rész, 18. oldal) csomópontra átvisszük, és az összegyűlt mennyiségeket a csomópontokban összegezzük. Ennek a módszernek hátrányos tulajdonsága, hogy azokban a pontokban, melyekből más pontok több úton is elérhetők, az összeggyakoriság többszöröse gyűlik össze.

Ennek elkerülésére követni kell, hogy egy pontba ugyanaz a gyakorisági érték több helyről fut-e be, és ha igen, akkor is csak egyszer szabad figyelembe venni. Ez utóbbi módszert nevezhetjük **alszerkezet**-örököltetésnek. A korpuszgyakoriság ( $f$ ) fogalmát kiterjeszthetjük az összes isz-re (azaz az szhb szerkezetekre is): ezeknél a korpusz minden olyan tagmondata bele fog számítani az isz gyakoriságába, amire az adott isz illeszkedik. Vegyük észre, hogy

az alszerkezet-örökltetés alkalmazásával az összes isz-re a korpuszgyakoriságát kapjuk meg (ld. a 15. ábra  $l = 4$  szintjét), azaz jelen esetben  $q \equiv f$ . Ez triviális művelet, ehhez nem is volna szükség a hálóra, csak meg kell számolni a korpuszban a kívánt szerkezeteket. Ekkor nyilván  $q(o_3) > q(o_1)$  és  $q(v_3) > q(v_1)$  lesz, valamint tudjuk, hogy  $o_3 \leq o_1$  és  $v_3 \leq v_1$ , azaz azt mondhatjuk, hogy „ha egy szerkezet (hálói értelemben) kisebb, akkor (gyakorisági mérőszám szerint) gyakoribb”.

Bevezetünk egy fontos elvet. A „gyakoriságmegmaradás elve” azt mondja ki, hogy kezdetben és az örökltetés során is a korpusz minden tagmondatát pontosan egyszer kell számolnunk.

Az alszerkezet-örökltetés esetén ez nem teljesül, hiszen minden mondat annyiszor számítottik, ahány alszerkezete van a mondatvázának (a 15. ábrán a *vet rám pillantást* 66-os gyakorisági értéke például a *vet rám vmit*-nél és a *vet vmire pillantást*-nál is megjelenik). Az elv teljesüléséhez egy fontos tulajdonsággal kell rendelkeznie az örökltetési eljárásnak, ún. **megjegyzős** örökltetésnek kell lennie. Ez két dolgot jelent. Egyrészt, hogy minden gyakorisági érték mindig csak egyfelé (1 konkrét 1-gyel kisebb csomópontra) öröklődhet tovább. (A nem-triviális (18. oldal) szerkezetek azok, ahol az örökltetés iránya kérdésként felmerül.) Másrészt mindig tudni kell, azonosíthatónak kell lennie, hogy adott eredeti korpuszmondathoz (vagy mondatvázhoz) rendelhető gyakorisági érték épp a háló melyik csomópontjánál van nyilván-tartva.

Azt, hogy az elv kezdetben is teljesüljön az biztosítja, hogy a korpuszgyakorisági értéket a tk szerkezetekhez rendeljük hozzá (6.4. rész). A gyakoriságmegmaradás elvének köszönhető, hogy a gyakorisági mérőszámok is ténylegesen gyakoriságok lesznek (a  $q$  értékek mindig konkrétan  $q$  darab eredeti korpuszmondathoz felelnek meg), csak nem az eredeti tk szerkezetek gyakoriságai, hanem a háló alkalmasan kiválasztott csomópontjainál lévő, az eredeti tk szerkezetekre illeszkedő szerkezetek (a majdani v-isz-ek!) gyakoriságai.

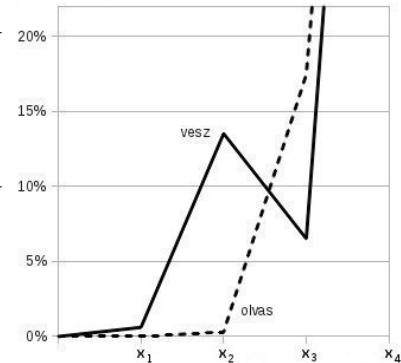
Az alszerkezet-öröklődés nem mondana semmit a v-isz-ekről, a szándékunk pedig éppen a v-isz-ek számba vétele és egyúttal gyakoriságuk megállapítása. A cél nem csupán a v-isz-ek feltérképezése általában, hanem lehetőleg az is, hogy minden egyes mondatvázra mondjuk meg, hogy ő konkrétan mely v-isz-nek a példánya. Ezért valamiféle megjegyzős öröklődést fogunk választani. A háló egyes csomópontjainál összegyűlt gyakoriság egyrészt ki fogja jelölni a v-isz-eket, másrészt a fentiek miatt a gyakoriságukat is rögtön meg fogja adni.

Fontos hangsúlyozni a korpuszgyakoriság ( $f$ ) és a gyakorisági mérőszám ( $q$ ) alapvető különbségét. A  $f$  az adott formális igei szerkezet előfordulásainak száma a korpuszban, a  $q$  viszont egy külön algoritmussal az igei szerkezetekhez rendelt érték, ami célja szerint a v-isz-eket ragadja meg, egész pontosan, hogy az adott isz hány esetben v-isz, azaz jellemző globális v-isz esetén magas, kevésbé jellemző globális v-isz esetén alacsony. Ebből következik, hogy itt már nem feltétlenül igaz, sőt általában nem igaz, hogy a kisebb szerkezet a (gyakorisági mérőszám szerint) gyakoribb, azaz  $x \leq y \not\Rightarrow q(x) > q(y)$ .

A mondatvázakhoz tehát v-isz-t kell rendelnünk. A következőt fogjuk tenni. A mondatvázakat valami módon a gyakoriságukkal együtt addig *visszük lefelé a hálón* – más szóval a gyakoriságokat örökltetjük, megjegyezve, hogy melyik mv-ból származnak –, amíg az eredeti mv v-isz-éhez nem érünk. A gyakoriságot *csak* ezen a ponton vesszük figyelembe, azt mondjuk, hogy a v-isz annyiszor fordul elő (annyi a gyakorisága), amennyi gyakoriság – persze akár számos mv-ból – összegyűlt nála. Ezen a módon a gyakoriságok a v-isz-eknél fognak felhalmozódni, mivel a hosszabb szerkezetektől az örökltetés során elveszük a gyakoriságokat, a v-isz-ek alatti rövidebb szerkezetekhez pedig eleve nem juttatunk. Ez tehát a v-isz-eket megállapító algoritmus kerete.

Vegyük észre, hogy a fenti *olvas*–*vesz* párhuzam a *v*-isz-ek keresése szempontjából megtévesztő, a hasonlóság csakis formai. Sőt, a két ige összevetése alapvető fontosságú különbségre mutat rá, és egyben megvilágítja a fent leírt algorímus-keret működését. Tekintsük át a 3. táblázatot.

igei szerkezet	jel	<i>f</i>	rel%	$\Delta\%$	<i>q</i>
<i>olvas könyvet ágyban</i>	<i>o</i> <sub>1</sub>	1	0.0%	0.0%	0
<i>olvas könyvet vmiben</i>	<i>o</i> <sub>2</sub>	78	0.3%	0.3%	0
<i>olvas vmit</i>	<i>o</i> <sub>3</sub>	4560	17.8%	17.5%	0
<i>olvas vmit</i>	<i>o</i> <sub>4</sub>	25653	100.0%	82.2%	25653
<i>vesz részt vitában</i>	<i>v</i> <sub>1</sub>	878	0.6%	0.6%	0
<i>vesz részt vmiben</i>	<i>v</i> <sub>2</sub>	20469	14.1%	<b>13.5%</b>	20469
<i>vesz vmit vmiben</i>	<i>v</i> <sub>3</sub>	29902	20.6%	6.5%	0
<i>vesz vmit</i>	<i>v</i> <sub>4</sub>	144812	100.0%	79.4%	124343



3. táblázat. Az *olvas* és a *vesz* egyes formailag egymásnak megfelelő szerkezeteinek összevetése. Mindkét igenél tárgyi és *-ban/-ben*-ragos hely szerepel, a 4 szerkezet rendre 5, 4, 3 és 2 hosszúságú, és a rövidebbek illeszkednek a hosszabbakra. A korpuszgyakoriság (*f*) esetében a lentebbi számok természetesen *tartalmazzák* a fentebbieket. A rel% oszlopban a relatív gyakoriság szerepel a tárgyias szerkezet százalékában. A  $\Delta\%$  oszlopban a megelőző és az aktuális rel% különbsége, azaz, hogy mennyit tesz hozzá az aktuális szerkezet az előzőekhez. A gyakorisági mérőszám (*q*) oszlopban az ideális értékek szerepelnek, itt a lentebbi számok természetesen *nem* tartalmazzák a fentebbieket,  $124343 = 144812 - 20469$ . A grafikonon *y* tengelyén a  $\Delta\%$ , *x* tengelyén pedig a 4-4 szerkezet szerepel. A grafikonon látható csúcsot a táblázatban vastagított számmal jelentjük meg.

Mit mond ez a táblázat? Azt, hogy a felírt szerkezetek közül az *olvas* esetében egy, a *vesz* esetében pedig kettő *v*-isz van. Az *olvas* szerkezeteit mind *olvas vmit*-ként számoljuk el, a *vesz*-nél azonban egy jelentős hányad esetében azt gondoljuk, hogy a *vesz részt vmiben* lesz a helyes *v*-isz, és csak a többi marad meg mint *vesz vmit*. Helyesen döntünk, ha így állapítjuk meg a *v*-isz-eket: a szóban forgó szerkezetek közül ez a három felel meg a *v*-isz definíciójának (24. oldal). Arról volt már szó, hogy a *vesz részt vmiben* nem alete, nem realizációja a *vesz vmit* szerkezetnek, ezért önállóan, külön kezelendő (25. oldal), további érveket ezzel kapcsolatban Sass (2011) 23-24. oldalán olvashatunk.

De honnan tudjuk, hogy hol kell megállni, hol kell befejezni a gyakoriságok öröklötését? Általánosságban érdemes törekedni arra, hogy minél magasabb szinten, azaz minél nagyobb *l* szerinti hosszánál álljunk meg, hogy ne alsóbb szinteken, rossz esetben a gyökérszintnél gyűljön össze az összes gyakoriság. Ezen felül úgy tűnik, hogy a grafikonon látható csúcs van segítségünkre. Amikor jelentős növekedés után jelentős visszaesés (azaz lokális maximum) van a  $\Delta\%$  grafikonjának menetében, akkor az egy *v*-isz-t jelezhet, érdemes lehet megállni. Azaz addig érdemes öröklötetni, amíg a (meglévőhöz képest) jelentős mennyiségű plusz gyakoriságot tudunk gyűjteni. Amikor nincs így, megállhatunk. Ez a tanulmány egyik legfontosabb megfigyelése.

A példa ezen pontja egyben azt a fenti állítást is szemlélteti, hogy nem mindig igaz, hogy a kisebb szerkezet a (*q* szerint) gyakoribb, mert  $v_3 \leq v_2$ , ugyanakkor  $q(v_3) \not\leq q(v_2)$ , tekintve, hogy az előbbi nulla, az utóbbi viszont több mint 20000. A *vesz vmit vmiben* szerkezetből a maradék  $29902 - 20469 = 9433$ -at persze öröklötjük tovább, ez megjelenik a *vesz vmit* szerkezet *q* értékében.

Ne felejtjük el, hogy a fenti példa is csak illusztrációnak tekinthető, mivel bár hiteles adatokat (Sass 2015) tartalmaz, a korpuszhálónak csak egy apró részletét mutatja be, a háló kontextusából kiragadva. Két darab 2-dimenziós 3-Post-hálóban haladtunk lefelé „cikkcakkban” (vö: 6. ábra, 23. oldal).

Két fontos kérdést hagytunk nyitva fent, mikor az algoritmus keretét vázoltuk: (1) pontosan mi legyen az öröklötési módszer, és (2) milyen feltétel esetén gondoljuk egy *v-isz*-ről, hogy *v-isz*, azaz megállhatunk az öröklötésben. E tényezők konkrét kidolgozásával lehet kialakítani a különféle egymással versengő, összevethető konkrét algoritmusokat.

Térjünk vissza a 35. oldalon látható 15. ábrához, és figyeljük meg a két felső szintet. Azt mondtuk, hogy kettő kívánt *v-isz*-t látunk: az  $l = 4$  szinten 67-es  $q$  értékkel szereplő *vet szemére vmit* szerkezetet és a 84-es  $q$  értékkel szereplő *vet vmire pillantást* szerkezetet. Az ábrán alszerkezet-öröklötés történik, ami megmutatja, hogy megjegyzős öröklötéssel maximum mennyi gyakoriság juthat az egyes csomópontokba. A nyilak vastagsága az öröklődő gyakorisági értékek nagyságát jelzi.

Ha megvizsgáljuk a *v-isz*-ek és a nyilak, illetve nyílvastagságok viszonyait, két szabályt állíthatunk fel. Azok a csomópontok lesznek jó eséllyel *v-isz*-ek, amelyek megfelelnek az alábbi feltételeknek.

#### **V-isz feltételek:**

- (1) nagy gyakorisági mérőszámmal bíró pontokból; és
- (2) sok pontból öröklötethető gyakoriság a csomópontba.

Azokról a csomópontokról van tehát szó, ahová sok vastag nyíl vezet. Azt látjuk, hogy éppen az említett két *v-isz* lesz az, ahová sok gyakoriság gyűlik „sok” (az ábrán lévő részlet kis mérete miatt csupán 2) helyről. Ha csak az egyik feltétel teljesül, nem *v-isz*-t kapunk: a *vet vmire ezt* ( $q = 49$ ) esetében csak (2), a *vet rám vmit* ( $q = 66$ ) esetében csak (1) teljesül. Ezekből a csomópontokból tovább öröklötethető a gyakoriság. Meg kell fontolni, hogy ezt a két feltételt hogyan tudjuk egyszerre figyelembe venni, és milyen súllyal vegyük őket figyelembe. Bizonyos esetekben önmagában elegendőnek tűnik az egyik feltétel megléte, ha az kellően erős: a *fest ördögöt falra mv* esetében – mivel mondatvázzról van szó – csak (1) teljesülhet, mégis érdemes elfogadni *v-isz*-ként, kiemelkedő gyakorisági értékével érvelve. Itt is megemlítjük, hogy az ábra csak egy kis méretű részhálót mutat be, és annak is csak két szintjét vizsgáltuk. A gyakoriságok valójában egymással kölcsönhatásban a teljes korpuszháló számos szintjén át öröklődnek. Elmondhatjuk, hogy a fentiek alapján valóban a két kívánt *v-isz* jön ki a 15. ábra példáján.

A 3. táblázat kiegészítő szemponttal szolgál a fenti alapvetéshez: a *v-isz* feltételek statikusan ragadják meg az esélyes csomópontokat, az *olvas-vesz* példájának tanulsága pedig az öröklötés folyamatában dinamikusan alkalmazható feltételt ad ahhoz, hogy adott ponton érdemes-e további öröklötési lépést végezni.

Kiemeljük, hogy a modell és az algoritmus, a tanulmányban megfogalmazott teljes rendszer az összes *v-isz*-t egységesen, egyformán kezeli. Nem diszkriminál a valódi igei szerkezetek között sem aszerint, hogy hány hely van bennük, sem aszerint, hogy van-e bennük bizonyos helyeken kitöltő vagy nincs. Ugyanolyan erővel tud egy vagy többemű entitáshoz (egyszerű vagy komplex ígéhez) további elemet (bővítményt, vonzatot) kapcsolni, ahogy ezt a példában is láttuk. Függetlenül attól is, hogy az adott esetraggal megjelölt hely egyszer a komplex ígé-

ben a kollokátumot jelöli ki, máskor meg a vonzati helyet. Az *olvas vmit* – ami egészen más struktúrájú, egyszerűbb szerkezet, mint a 15. ábra  $l = 4$  szintjén megszülető *vet szemére vmit* és *vet vmire pillantást* – pontosan ugyanazon a módon fog kijönni, mint az ábrán lévők. Az *olvas* ige hálójában  $l = 3$  szinten is még annyiféle különböző elem lesz (*könyv, újság, vers, cikk, regény, írás, szöveg, mű, biblia, levél ...*), hogy a *v-isz* feltételek – főként a (2) – miatt érdemes továbbörököltetni a gyakoriságot  $l = 2$  szintre, ahol megkapjuk az *olvas vmit v-isz-t*.

A számítógépes nyelvészetben, nyelvtechnológiában régi hagyománya van az „egyszavas” és a „többszavas” egységek különválasztott kezelésének. Most amellet érvelünk, hogy bizonyos esetekben érdemes elvonatkoztatni attól, hogy valami hány szóból áll, és a kitűzött feladat megoldása során elfogadni, hogy az eredmény néha egyszavas, máskor pedig többszavas. A modellünk egyik fő gondolata, hogy ne válasszuk szét csupán formai alapon az egyszavas és a többszavas konstrukciókat, hanem kezeljük egységesen mindet. Ne mondjuk meg előre, hogy egyszavas, többszavas vagy hány elemből álló egységeket várunk, hanem nézzük meg, hogy a korpuszban mi van. Esetünkben az a cél, hogy rájövünk, hogy egy mondatban mi a *v-isz*; sokdrangú kérdés, hogy az éppen hány szóból vagy elemből épül fel. Mindegy, hogy egy *v-isz* vonzatos egyszerű ige, vagy éppen vonzatos komplex ige, meghatározása pontosan ugyanolyan fontos, nem tűnik észszerűnek a kettőt külön feladatként megközelíteni. Ráadásul a szavak száma csak egy felszíni jellemző, valójában a helyek és kitöltők elrendeződése számít.

Az egyik üzenet az lehet, hogy az összetett egységek nem szavakból, hanem helyekből és kitöltőkből állnak, hogy ezt a két típust érdemes különválasztani (vö: 5.3. rész, 25. oldal). Esetleges, hogy egy *v-isz* egy nyelvben épp többszavas kifejezés vagy nem. A *participate in* vonzatos egyszerű ige, a *részt vesz vmiben* vonzatos komplex ige, többszavas kifejezés. A *participate in debate* *isz* és a *vesz részt vitában* *isz* véletlenül egyaránt három szó. De ha jobban megnézzük a modellünkben, akkor az angol szerkezet *1.ige + 2.hely + 3.kitöltő* struktúrájú, a magyar pedig ettől teljesen eltérő *1.ige + 2.kitöltő-hely + 3.kitöltő-hely* struktúrájú úgy, hogy ez utóbbiban az első két szó komplex igét alkot. Így is mondhatjuk: alapegységünk a szerkezet, azaz a több elemből álló egység. Ez különleges esetként magában foglalja azt is, ha néha egy szerkezet csak egy elemből áll.

A korpuszháló éppen olyan struktúra, amit egy *v-isz*-eket kereső algoritmusnak érdemes alapul vennie. Ennek az a gyökere, hogy a *v-isz*-ek kitöltött és a kitöltetlen helyeit (vagy ha úgy tetszik a kollokációkat és a vonzatokat) pontosan ugyanazokkal a nyelvi elemekkel, formai eszközökkel (esetragokkal, névutókkal stb.) fejezzük ki, azaz a felszínen ezek ugyanúgy néznek ki. Azaz sose tudható előre, hogy melyik hely lesz végül kitöltött és melyik nem, erről az algoritmusnak kell döntenie, konkrétan hogy továbbörököltet egy kitöltőről, és ezáltal elhagyja, vagy nem.

Említettük (36. oldal), hogy a hálóalapú *v-isz*-felfedező eljárásunkkal a globálisan jellemző *v-isz*-eket fogjuk megkapni, nem a konkrét mondatban megjelenő lokális *v-isz*-t. Ha jobban meggondoljuk, a megjegyzős örököltetés során mégiscsak meghatározzuk a lokális *v-isz*-t is. A megfelelő mondatvázak gyakoriságait összesítve kapjuk meg a *v-isz*-ekhez tartozó gyakorisági mérőszámot. Ugyanakkor – éppen a megjegyzésnek köszönhetően – ismerjük, hogy az adott *v-isz*-nél összegyűlt gyakoriság mely mondatvázakból, végső soron mely mondatokból származik. Azaz van információnk arról, hogy egy adott mondathoz milyen *v-isz* tartozik: az, amelynél az adott mondat a *v-isz* gyakorisági mérőszámába bele van számítva. Ezen a módon a megjegyzős örököltetési eljárások egyszerre megoldják a globális és a lokális *v-isz*-azonosítási feladatot.

Említettük (5.2. rész), hogy a globális  $v$ -isz-ek azok az igei szerkezetek, melyeket érdemes lehet egy szótárban, illetve kifejezéstárban összegyűjteni. A munkálatok egyik végső eredménye éppen ez lehet: egy nyelv valódi igei szerkezeteit tartalmazó szótár.

### 7.3 Egy korábbi megközelítés

Egy korábbi dolgozatban (Sass 2011: 3.3. rész) szerepel egy  $v$ -isz-ekhez hasonló konstrukciók kinyerését célzó algoritmus. Nem illik bele maradéktalanul a jelen tanulmányban vázolt hálól alapú keretbe, de összevetni érdemes vele, mert sok közös pont, hasonló gondolat van. Az említett dolgozatban leírtak sok szempontból esetlegesnek, illetve barkácsolásnak tűnnek. Hiányzik a szilárd formális megalapozás, amit éppen a jelen tanulmány kíván kidolgozni és megadni, és egyúttal továbbgondolni és általánosítani.

Az első szembeötlő különbség, hogy nincs meg a hálófogalom, az igei szerkezetek hálós modellje, következésképpen nincs felépítve a korpuszháló sem. A módszer a mondatvázak listáját csak néhány hasznosnak tűnő szhb szerkezettel egészíti ki, azaz egy hiányos, „lyukacsos” hálót használ. Egész pontosan a mondatvázakon és az azokból képzett tklen isz-eken kívül a 2-helyes mondatvázakhoz (kizárólag a 2-helyesekhez!) hozzáveszi az  $l = 4$  szintű, szabad helyel bíró szerkezeteket is (vö: 6. ábra, 23. oldal) az általa opcionalizálásnak nevezett művelettel (ti. „opcionálissá” teszi a kitöltőket, így jönnek létre ezek a szerkezetek). Ennek köszönhetően esélyt ad arra, hogy a kiemelten fontos *vesz részt vmiben* típusú komplex  $v$ -isz-ek is megjelenhessenek az eredményben.

Nem is a hálón lépeget az algoritmus, hanem a mondatvázak (szhb és tklen szerkezetekkel a fentiek szerint kiegészített) listáján, a mondatvázakat veti össze rendre egymással, így a mondatvázak számában (ami akár milliós is lehet) négyzetes működésű, ami nagyon nem hatékony. (Ezen segíthet, hogy a hálóban viszont mindig minden szerkezet közvetlenül össze van kapcsolva a nála 1-gyel rövidebb szerkezetekkel, azok közvetlenül elérhetők belőle.)

Az öröklötési módszerét nevezzük **véletlenszerű** öröklötésnek: minden mondatvázból a gyakoriságot az illeszkedő, rövidebb (a rövidebbek között a lehető leghosszabb) szerkezetek közül egy véletlenszerűen kiválasztott kapja. (Kiemelendő, hogy a véletlenszerű öröklötés megjegyzős.) A módszer abban bíz, hogy a jó helyeken gyűlik össze sok gyakoriság. Ezt arra alapozza, hogy szabad helyeken általában sokféle kitöltő szerepel, és ezért sok irányból jöhet a gyakoriság egy ilyen szabad helyhez (vö: (2)-es  $v$ -isz feltétel a 39. oldalon). A véletlenszerű öröklötés miatt szükség volt egy visszaellenőrzésnek nevezett trükkre, ami mondatvázanként a túl rövid szerkezetekhez tévedt gyakoriságokat hivatott utólag a helyükre terelni, a hosszabb, specifikusabb szerkezetekhez. Az eljárás nagyon egyszerű megállási feltételt alkalmaz: ha az összegyűlt gyakoriság a fix 5 értéknél több, akkor megállunk, nem öröklötünk tovább, a szerkezetet elfogadjuk  $v$ -isz-nek. Ez a jelen tanulmányban leírtak szerint nem tűnik jó ötletnek, tekintve, hogy a  $v$ -isz feltételek mindkét pontját megszegi sok esetben, mégis elfogadható végső eredményhez vezet.

Egyéb trükköket is alkalmaz a módszer. Az alanyi szabad helyet nem tartja nyilván, törli, más oldalról nézve úgy is mondhatjuk, hogy mindig odaérti. Ez ad lehetőséget arra, hogy az említett 2-helyes,  $l = 4$ , szhb kategóriába sok fontos szerkezet beférjen. A magyar nyelv sajátosságaiából adódóan valójában nem csak tk isz-kre van gyakorisági adatunk, hanem az olyan egy szabad tárgyi hellyel bíró isz-kre is, ahol a tárgy léte csak az ige határozott ragozása utal. Ezeket számításba veszi a módszer. Kitöltőként gyakran megjelennek névmások. Mivel ezek általában nem  $v$ -isz részei, ezért ezeket utólag („kézzel”) törlik, és összevonják a törlés

után azonossá váló *isz*-ket, a megjegyzősséget megtartva a gyakorisági mérőszámok egyszerű összegzésével.

Ez a módszer megszegi a teljesség (vö: 24. oldal) követelményét is, mivel a háló lyukacsos-sága miatt bizonyos, az opcionalizálás által meg nem engedett hosszabb szerkezeteket eleve kizár, így ezek elvesznek. Ilyenek például a 3-helyes vegyes szerkezetek, köztük a *vesz őri-zetbe vkit vmi-miatt*. A teljesség feladásáért cserébe a futási idő kezelhetőbb marad. (Ezen is segít a hálós megközelítés, segítségével ragaszkodhatunk a teljességhez.)

Érdekes módon ez az esetlegesnek tűnő módszer is viszonylag jó eredményeket ad. Bár sok esetben megjelennek a vélelmezett *v-isz*-ekben szabad határozók, gyakori esetragok (pl. *-bAn*), illetve az esetleges kitöltők, mégis általában és elsősorban – ahogy a kiértékelés mutatja – *v-isz*-ek jönnek ki, pontosabban megkapjuk a „lehető legspecifikusabb ugyanakkor elegendően gyakori” *isz*-ket (Sass 2011: 60).

A módszerrel egy kísérleti korpuszvezérelt szótár (Sass et al. 2010) is készült, ami már mutatja azokat a gyakorisági mérőszámok között megjelenő viszonyokat, amire a 3. táblázatban (38. oldal) utaltunk. A várt pontokon rendre előfordul, hogy egy hosszabb szerkezet (ami ténylegesen globális *v-isz*) nagyobb gyakorisági mérőszámmal bír, mint egy rövidebb (amit kevésbé, vagy egyáltalán nem gondolnánk *v-isz*-nek):

$$q_{sz}(v_3 : \text{vesz vmit vmiben}) = 840 < q_{sz}(v_2 : \text{vesz részt vmiben}) = 3985$$

ahol  $q_{sz}$  a szótárban lévő gyakorisági mérőszám. A módszer korlátaiból adódóan  $q_{sz}(v_3) \neq 0$ , de a gyakorisági viszonyok megfelelőek.

Látjuk, hogy a most bemutatott módszer egyszerűsítéseket, esetleges megoldásokat alkalmaz, bizonyos megfontolásokat nem vesz tekintetbe, az eredmény mégis egészen jó (Sass 2011: 3.3.2. rész). Ez a módszer csak egy gyors, közelítő megoldásnak tekinthető. Mostantól, a hálós modell felhasználásával, lehetőség van megalapozottabb és jobb algoritmusokat tervezni, készíteni a valódi igei szerkezetek feltérképezésére.

## 8 Összefoglalás

A tanulmányban bemutattunk egy absztrakt algebrai modellt, az ún. duplakocka modellt. Ez egy olyan hálóstruktúra, amely megfelelő számú egymásra épülő  $n$  dimenziós kockaként kezelhető el (2–4. rész). Megállapítottuk, hogy a bevezetett objektumok az ún.  $h$  dimenziós ( $h = 1, 2, \dots$ ) harmadrendű Post-hálókkal izomorfak.

Ezt az absztrakt modellt egy bizonyos módon konkretizáltuk: igei szerkezetek (egy tagmondat igeje és a mellette álló bővítmények, kitöltetlen és kitöltött helyek) ábrázolására használtuk. Megmutattuk, hogy éppen ez az a struktúra, ami megfelel a célnak. Így kaptuk az igeiszerkezet-hálókat (5. rész). Az igeiszerkezet-hálók egy különleges pontját megjelölve bevezettük a valódi igei szerkezet fogalmát, bemutattuk ennek jelentőségét egy nyelv szótára szempontjából.

Kialakítottunk egy módszert, amelynek segítségével nem csak egy tagmondatot (egy „igés egységet”), hanem egy teljes szöveget lehet a modell segítségével ábrázolni. Az igeiszerkezet-hálóból mint alapelemekből felépítettük az ún. korpuszhálót (6. rész). A korpuszháló felépítése kapcsán leírtunk egy izgalmas struktúrát, amelyben különleges prímszerű elemek fordulnak elő, melyeket magasabb dimenziós prímeknek neveztünk.

Végül, felismerve, hogy a valódi igei szerkezetek a korpuszháló bizonyos kitüntetett pontjain helyezkednek el, felvázoltuk a valódi igei szerkezetek felfedezésére szolgáló algoritmus

keretét, elemeit (7. rész). Megfogalmaztuk azt a sejtést, hogy a valódi igei szerkezetek a gyakorisági grafikon lokális maximumainál találhatók.

Ahogy utaltunk rá, a tanulmányban az absztrakt modellünket – annak bemutatása után – *egyetlen* konkrét esetre alkalmaztuk: magyar igei szerkezetek kezelésére. Az absztrakt modellben benne van a lehetőség, hogy az iménti kijelentést annak minden pontján általánosítsuk. A modell a vázolt algoritmussal együtt különösebb bonyodalmak nélkül alkalmazható egyéb nyelvekre. Lényegében csak annyi szükséges, hogy az igék melletti helyeket kell valahogyan karakterizálnunk, ahol a különféle bővítmények elhelyezkednek. A magyar eseteken és névutókon túl helyeket megadhatnak előljárók, előljáró–esetrag kombinációk, vagy például szórendi megkötések is. Az esetragos/névutós névszói bővítményeken kívül egyéb elemek (például háttározószó, igekötő) helyként való számításba vételét szintén meg lehet fontolni. Sőt, a modellt nem-igei (hanem például főnévi) központú szerkezetek kezelésére is lehet használni, illetve lényegében bármilyen olyan – akár nem nyelvi – struktúra kezelésére, ami valamilyen módon megfeleltethető a komód-hasonlatnak (13. oldal), és jelentős mennyiségű adat („korpusz”) áll rendelkezésre belőle.

A duplakocka modell egyszerű és általános: matematikailag jól megfogható, ugyanakkor mindenféle igei szerkezetet magában foglal, egységesen kezel. Azon túl, hogy új szemléletet ad az igei konstrukciók vonatkozásában, alapot adhat a valódi igei szerkezetek automatikus felfedezéséhez és hozzájárulhat az igei szerkezetek jellemzőinek, természetének, kompozicionálisának–idiómaságának kutatásához.

### Köszönetnyilvánítás

A kutatást az MTA Bolyai János Kutatási Ösztöndíja támogatta (ügyszám: BO/00064/17/1; időtartam: 2017-2020). Köszönet Sain Ildikónak és Sági Gábornak a duplakocka modellről, általában a hálókról és konkrétan a Post-hálókról szóló beszélgetésért.

### Irodalom

- Atkins, B.T.S. & Rundell, M. (2008): *The Oxford Guide to Practical Lexicography*. Oxford: Oxford University Press.
- Biedermann, K. (1998): Powerset trilattices. In: Mugnier, M.-L. & Chein, M. (eds.): *Conceptual Structures: Theory, Tools and Applications*. Heidelberg: Springer, 209-244.
- Čech, R., Pajas, P. & Mačutek, J. (2010): Full valency. Verb valency without distinguishing complements and adjuncts. *Journal of Quantitative Linguistics*, 17(4): 291-302.
- Chechik, M., Devereux, B., Easterbrook, S., Lai, A.Y.C. & Petrovykh, V. (2001): Efficient multiple-valued model-checking using lattice representations. *Proceedings of CONCUR'01*.
- Epstein, G. (1960): The lattice theory of Post algebras. *Transactions of the American Mathematical Society*, 95(2): 300-317.
- Epstein, G. (1993): *Multiple-Valued Logic Design: an Introduction*. Bristol: IOP Publishing.
- Kalmár, L. (1964): Matematikai és nyelvi struktúrák. *Általános Nyelvészeti Tanulmányok*, 2: 11-74.
- Kovács, L., Orosz, K. & Pollner, P. (2012): Magyar szóasszociációk hálózata. *Magyar Tudomány*, 173(6): 699-705.
- Mac Lane, S. (1998): *Categories for the Working Mathematician*. Heidelberg: Springer.



- Mazur, B. (2008): When is one thing equal to some other thing? In: Gold, B. & Simons, R.A. (eds.): *Proof and Other Dilemmas: Mathematics and Philosophy*. Series: Spectrum, 221-241.
- Miháltz, M., Hatvani, Cs., Kuti, J., Szarvas, Gy., Csirik, J., Prószéky, G. & Váradi, T. (2008): Methods and results of the Hungarian WordNet project. In: *Proceedings of The Fourth Global WordNet Conference*. Szeged, 311-321.
- OEIS Foundation Inc. (2011): *The On-Line Encyclopedia of Integer Sequences*. <http://oeis.org>.
- Pagliani, P. & Chakraborty, M. (2008): *A Geometry of Approximation: Rough Set Theory: Logic, Algebra and Topology of Conceptual Patterns*. Heidelberg: Springer.
- Partee, B.H., Ter Meulen, A. & Wall, R.E. (1990): *Mathematical Methods in Linguistics*. Dordrecht, Boston & London: Kluwer Academic Publishers.
- Sag, I., Baldwin, T., Bond, F., Copestake, A. & Flickinger, D. (2002): Multiword expressions: A pain in the neck for NLP. In: *Proceedings of 3rd CICLING*, 1-15, Mexico City, Mexico.
- Sass, B. (2011): *Igei szerkezetek gyakorisági szótára – egy automatikus lexikai kinyerő eljárás és alkalmazása*. PhD disszertáció. PPKE ITK.
- Sass, B. (2015): 28 millió szintaktikailag elemzett mondat és 500000 igei szerkezet. In: *XI. Magyar Számítógépes Nyelvészeti Konferencia (MSZNY2015)*. Szeged: JATEPress, 303-308.
- Sass, B., Váradi, T., Pajzs, J. & Kiss, M. (2010): *Magyar igei szerkezetek – A leggyakoribb vonzatok és szókapcsolatok szótára*. Budapest: Tinta Könyvkiadó.
- Shramko, Y., Dunn, J.M. & Takenaka, T. (2001): The trilattice of constructive truth values. *Journal of Logic and Computation*, 11(6): 761-788.
- Tesnière, L. (2015): *Elements of Structural Syntax*. Amsterdam: John Benjamins.
- Teubert, W. (2005): My version of corpus linguistics. *International Journal of Corpus Linguistics*, 10(1): 1-13.
- Tognini-Bonelli, E. (2001): *Corpus Linguistics at Work*. Amsterdam: John Benjamins.
- Vincze, V. (2011): *Semi-Compositional Noun + Verb Constructions: Theoretical Questions and Computational Linguistic Analyses*. PhD disszertáció. Szegedi Tudományegyetem.

Sass Bálint PhD  
MTA Nyelvtudományi Intézet  
Nyelvtechnológiai Kutatócsoport  
H-1394 Budapest  
Pf. 360  
[sass.balint@nytud.mta.hu](mailto:sass.balint@nytud.mta.hu)