



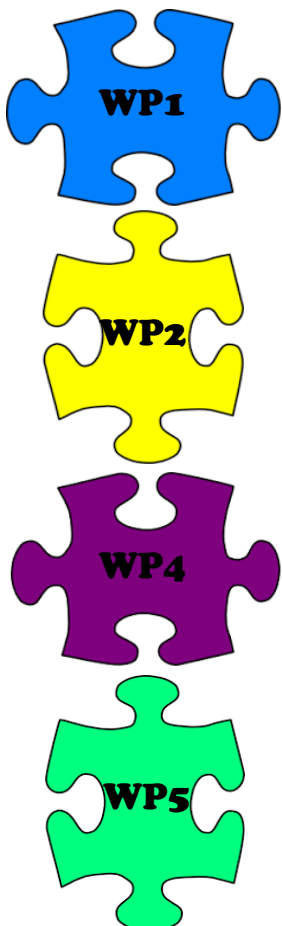
### Cyfarchion gan y Prif Ymchwilydd

*Mae trydydd mis prosiect CorCenCC yn mynd rhagddo'n dda, ac mae'n braf gallu dweud bod y gwaith yn llwyddo. Yr ydym yn dechrau gwneud cynnydd sylweddol ar becynnau gwaith unigol (WPs) sy'n rhan o'r prosiect, ac yr wyf yn falch o ddweud ein bod eisoes wedi cyflawni llawer yn y cyfnod cymharol fyr hwn. Yng nghylchlythyr y mis hwn byddwn yn darparu diweddariadau byr ar rai o'r datblygiadau sy'n parhau ar draws Pecynnau Gwaith penodol yn y prosiect, a byddwn yn arddangos yn arbennig waith diweddar ar ddatblygiad offeryn tagio semantig ar gyfer y Gymraeg sy'n cael ei arwain gan y tîm ym Mhrifysgol Caerhirfryn (Lancaster) (WP3). Yn ogystal â hyn, y mis hwn cynhwysir colofn newydd a fydd yn eich cyflwyno i aelodau unigol tîm prosiect CorCenCC.*

*Mwynhewch! Dr Dawn Knight (Prifysgol Caerdydd)*

### Diweddariadau ar Becynnau Gwaith (WPs) CorCenCC 1, 2, 4, 5

Mae'r gwaith ar brosiect CorCenCC wedi'i ddsbarthu ar draws 6 phecyn gwaith cydgyssylltiedig, y mae gan bob un ohonynt dasgau, nodau ac amcanion penodol. Mae WPO yn cynnwys gweithgareddau dylunio, cwmpasu a gweithgareddau hyfforddi, ac yn cynnwys holl aelodau tîm y prosiect. Ceir diweddariadau byr ar Becynnau Gwaith 1, 2, 3, 4, isod (gweler yr adran nesaf ar gyfer diweddariadau ar WP3).



#### **Nod: Casglu, trawsgrifio a dileu manylion enwau o'r data (arweinydd: Steve Morris)**

Mae tîm WP1 wedi bod yn datblygu'r fframwaith samplu ar gyfer data sydd i'w gasglu ar gyfer CorCenCC. Yn dilyn ymchwil helaeth, mae fersiwn ddrafft o hyn wedi'i llunio a'i dosbarthu i arbenigwyr ar y Gymraeg a'r corpws am adborth. Mae'r tîm hefyd wedi bod yn gweithio ar lunio confensiynau trawsgrifio a dileu manylion enwau ar gyfer y corpws.

#### **Nod: Datblygu'r set tagiau/tagiwr rhannau ymadrodd (arweinydd: Dawn Knight)**

Mae cynnydd yn parhau ar ein set tagiau unigryw ar gyfer rhannau ymadrodd a'n hoffer tagio ar gyfer y Gymraeg. Ar ben hynny, rydym yn rhoi cynlluniau ar waith ar gyfer cynhyrchu set ddata safon aur yn ystod y misoedd nesaf, ar gyfer hyfforddi/gwerthuso offer prosesu ar gyfer Cymraeg naturiol!

#### **Nod: Cwmpasu/llunio'r pecyn cymorth pedagogaidd ar-lein (arweinyddion: Enlli Thomas/Tess Fitzpatrick)**

Mae gwaith wedi dechrau ar ystyried y mathau o offer ar-lein sydd eisoes ar gael i ddysgwyr Cymraeg er mwyn osgoi dyblygu adnoddau presennol, ac mae cynlluniau ar y gweill i gynnal arolwg i ymchwilio i ba adnoddau mae dysgwyr ac athrawon eisoes yn eu defnyddio a beth fydden nhw'n ddefnyddol yn hoffi ei weld yn cael ei ddatblygu fel rhan o becyn cymorth pedagogaidd CorCenCC.

#### **Nod: Adeiladu seilwaith i gynnal CorCenCC a chreu'r corpws (arweinydd: Irena Spasic)**

Mae ein gweinydd prosiect bron yn weithredol ac mae gwaith yn cychwyn i roi ein cymhwysiad torfoli data ar waith, sef y camau mawr cyntaf wrth osod seilwaith CorCenCC yn ei le (cadwch lygad ar [www.corcenc.org](http://www.corcenc.org) i gael gwybod am unrhyw ddatblygiadau).

## Datblygu Offeryn Anodi Semantig ar gyfer Dadansoddi'r Iaith Gymraeg

Ym Mhrosiect CorCenCC, sy'n ceisio llunio corpws mawr Cymraeg a datblygu offer ar gyfer yr iaith Gymraeg, rydym yn datblygu cyfres o feddalwedd prosesu ar gyfer y Gymraeg i gynorthwyo wrth ddadansoddi a chwilio am wybodaeth amrywiol a fydd wedi'i storio yn nata'r corpws.



Un o'r prif offer sy'n cael eu datblygu yn y prosiect hwn ym Mhrifysgol Caerhirfryn yw'r meddalwedd anodi semantig, sydd â'r nod o helpu i ddadansoddi data yn y Gymraeg ar lefel semantig, a hynny ar raddfa fawr. Bydd yr offeryn hwn yn prosesu samplau Cymraeg yn awtomatig ac yn labelu pob gair neu ymadrodd â'r ystyr bras. Bydd yn seiliedig ar system USAS Caerhirfryn (<http://ucrel.lancs.ac.uk/usas/>), ac yn ei hestyn, ac felly yn gallu pennu categorïau semantig ar ffurf tagiau megis I1.3 (Arian: Prisiau/Money: Price), K5.1 (Chwaraeon/Sports) i eiriau ac idiomau Cymraeg ac ymadroddion penodedig eraill (a elwir hefyd yn ymadroddion amleiriog) ar sail cynllun dosbarthiad semantig. Mae'r cynllun hwn yn cynnwys 232 o gategorïau semantig, sy'n syrthio i 21 o brif gategorïau, fel y dangosir isod:

Tag	Diffiniad	Tag	Diffiniad
A	TERMAU CYFFREDINOL A HANIAETHOL (GENERAL AND ABSTRACT TERMS)	N	RHIFAU A MESUR (NUMBERS AND MEASUREMENT)
B	B Y CORFF A'R UNIGOLYN (THE BODY & THE INDIVIDUAL)	O	SYLWEDDAU, DEFNYDDIAU, GWRTHRYCHAU AC OFFER (SUBSTANCES, MATERIALS, OBJECTS AND EQUIPMENT)
C	CELFF A CHREFFT (ARTS AND CRAFTS)	P	ADDYSG (EDUCATION)
E	GWEITHREDIADAU, CYFLYRAU A PHROSES AU EMOSIYNOL (EMOTIONAL ACTIONS, STATES & PROCESSES)	Q	GWEITHREDIADAU, CYFLYRAU A PHROSES AU IEITHYDDOL (LINGUISTIC ACTIONS, STATES AND PROCESSES)
F	BWYD A FFERMIO (FOOD & FARMING)	S	GWEITHREDIADAU, CYFLYRAU A PHROSES AU CYMDEITHASOL (SOCIAL ACTIONS, STATES AND PROCESSES)
G	LLYWODRAETH A'R CYHOEDD (GOVERNMENT AND THE PUBLIC DOMAIN)	T	AMSER (TIME)
H	PENSAERNIAETH, ADEILADAU, TAI A'R CARTREF (ARCHITECTURE, BUILDINGS, HOUSES & THE HOME)	W	Y BYD A'N HAMGYLCHEDD (THE WORLD AND OUR ENVIRONMENT)
I	ARIAN A MASNACH (MONEY & COMMERCE)	X	GWEITHREDIADAU, CYFLYRAU A PHROSES AU SEICOLEGOL (PSYCHOLOGICAL ACTIONS, STATES AND PROCESSES)
K	ADLONIAN, CHWARAEON A GEMAU (ENTERTAINMENT, SPORTS AND GAMES)	Y	GWYDDONIAETH A THECHNOLEG (SCIENCE AND TECHNOLOGY)
L	BYWYD A PHETHAU BYW (LIFE AND LIVING THINGS)	Z	ENWAU A GEIRIAU GRAMADEGOL (NAMES AND GRAMMATICAL WORDS)
M	SYMUD, LLEOLIAD, TEITHIO A CHLUDIANT (MOVEMENT, LOCATION, TRAVEL AND TRANSPORT)		

Mae datblygu offeryn anodi semantig cywir yn dasg heriol dros ben. Bydd ein hofferyn yn cyfuno sylfaen o wybodaeth eiriadurol semantig Gymraeg (yn ei hanfod geiriadur mawr y gall peiriant ei ddarllen) ac amrywiaeth o ddulliau **dileu amwysedd yn ystyr geiriau**, er mwyn cyflawni lefel uchel o gywirdeb yn yr anodi semantig. At hynny, byddwn yn archwilio dulliau gweithredu sy'n defnyddio torfoli, lle gwahoddir cymunedau Cymraeg eu hiaith i gymryd rhan, er mwyn helpu i greu'r sylfaen o wybodaeth eiriadurol semantig ar raddfa fawr a gwella perfformiad yr offeryn.

Gall yr offeryn anodi semantig ar gyfer y Gymraeg fod yn ddefnyddiol at amrywiaeth o ddibenion academiaidd ac ymarferol. Yn gyntaf, bydd y corpws Cymraeg cyfan sydd i'w greu ym Mhrosiect CorCenCC yn cael ei dagio gan ddefnyddio'r offeryn hwn, a fydd yn hwyluso amrywiol ddadansoddiadau semantig o ddata'r corpws a dull dibynadwy, cyflym o chwilio drwy ddata'r corpws yn ôl amrywiol categorïau semantig, fel y dangosir ar safle ymchwil corpws Wmatrix Caerhirfryn (<http://ucrel.lancs.ac.uk/wmatrix/>). Ar ben hynny, gellir defnyddio'r offeryn anodi semantig i wella systemau TGCh sy'n seiliedig ar yr iaith Gymraeg. Er enghraifft, gall ganiatáu i system TGCh echdynnu a dadansoddi gwybodaeth semantig benodol yn gyflym o ffynonellau Cymraeg ar-lein, a hynny ar raddfa fawr, gan helpu i gysylltu ffynonellau ar-lein sy'n ddefnyddiol ar gyfer defnyddwyr.



Er mai yn ddiweddar y cychwynnodd Prosiect CorCenCC, mae'r gwaith i ddatblygu'r offeryn anodi semantig yn mynd rhagddo'n dda, ac rydym wedi rhyddhau'r gyfres o categorïau wedi'i chyfieithu a welir uchod.

***Dr Paul Rayson (arweinydd WP3) a Dr Scott Piao (Cynorthwy-ydd Ymchwil), Prifysgol Caerhirfryn***

### **Cwrdd â'r tîm**

*Bob mis byddwn yn rhoi sylw i aelod gwahanol o'r tîm CorCenCC estynedig yn ein cylchlythyr. Bydd hyn yn rhoi cyfle i bawb ddweud ychydig wrthyfych chi am eu cefndir; beth maen nhw am ei weld o CorCenCC a sut maen nhw'n credu y gallai gyfrannu at eu gwaith eu hunain, neu yn fwy cyffredinol, at waith eraill yng Nghymru. Y mis yma mae'r sbotolau ar Dr Emyr Davies o CBAC-WJEC, sy'n rhan o Grŵp Ymgynghorol Prosiect CorCenCC.*

### **Proffil: Dr Emyr Davies**

O'r diwedd! Dyna fy ymateb i i'r newyddion y byddai gennym gorpws cynhwysfawr o'r iaith Gymraeg. Mae mwy nag un corpws bach wedi ei ddatblygu yn y gorffennol, ond dim byd ar y raddfa hon gyda chymaint o bosibiliadau iddo. Rwy'n ymddiddori mewn ieithoedd ac ieithyddiaeth ers amser maith. Ar ôl graddio yn y Gymraeg o Brifysgol Cymru, Aberystwyth ac ennill doethuriaeth, treuliais 11 mlynedd yn dysgu Cymraeg i oedolion yng Ngholeg y Drindod, Caerfyrddin. Treuliais lawer o amser yn datblygu adnoddau newydd, ond roedd hyn cyn y chwyldro digidol; doedd syniadau fel corpws electronig ddim yn bodoli. Cymerais flwyddyn i ffwrdd ym 1996 i astudio ar gyfer Uwch Ddiploma mewn Ieithyddiaeth yn UCD Dulyn, a mwynhau hwn yn fawr. Mae cwrs fel hwn yn rhoi sail gadarn i'r sawl sy'n gweithio ym maes ieithoedd, hyd yn oed os ydyn nhw'n mynd ymlaen i weithio mewn rhyw agwedd ar ieithyddiaeth gymhwysol yn ddiweddarach – a dyna'n union wnes i.



Gadewais i Gaerfyrddin i weithio yn CBAC, y bwrdd arholi â'i bencadlys yng Nghaerdydd, yn 2001, ac rwy'n gweithio'n bennaf ym maes asesu a datblygu adnoddau i oedolion sy'n dysgu Cymraeg ers hynny. Fy mhrif orchwyl oedd sefydlu cyfres o arholiadau iaith i'r sector, sydd bellach yn eu lle. Rhan arall o'm gwaith oedd cynrychioli'r arholiadau Cymraeg yn ALTE, Cymdeithas i Brofwyr Ieithoedd yn Ewrop. Mae hyn wedi bod yn brofiad gwych i mi, ac wedi rhoi cyfleoedd i fynd i weithdai gan rai o arbenigwyr blaenaf y maes ac i ddysgu ganddynt. Does dim llawer o bobl yn gweithio ar asesu yn unig, yn sicr yn Gymraeg, felly mae ALTE (yn un peth) yn cynnig cymuned broffesiynol o sefydliadau sy'n canolbwyntio ar asesu ieithoedd. Rhaid hefyd i aelodau ALTE gael awdit, neu

gyflwyno tystiolaeth ein bod yn cyrraedd proffil ansawdd uchel. Yr un safonau trylwyr a ddefnyddir ar draws pob iaith a phob sefydliad sydd eisiau bod yn aelod llawn, gan gynnwys pob prif iaith Ewropeaidd a llawer o'r ieithoedd llai eu defnydd, gan gynnwys y Fasgeg, yr Wyddeleg a'r Gymraeg.

Gweithio ar sail tystiolaeth yw cyrchddull ALTE. Mewn geiriau eraill, allwn ni gynnig tystiolaeth i gefnogi'r penderfyniadau a wnawn? Gallai corpws cynhwysfawr fod yn ddefnyddiol wrth ddatblygu profion, yn ogystal ag adnoddau ar gyfer dysgu, a byddai'n rhoi tystiolaeth i ni wneud mil a mwy o benderfyniadau. Daw llawer o gwestiynau i'r meddwl: *Sut mae dewis agweddau ar iaith i'w dysgu (a'u profi) gyntaf? Ydy tablau amllder o eirfa'r Gymraeg yn cyfateb i'r Saesneg? Oes cyfleoliadau y dylen ni eu cynnwys yn ein deunyddiau dysgu? Sut gallwn ni gynnwys iaith 'wastraff' mewn ffordd ddilys, e.e. wel...ym...? Pa agweddau ar dafodiaith sydd fwyaf defnyddiol i ddysgwyr? Ydy'r ymadrodd hwn a hwn yn ddefnyddiol mewn gwirionedd, neu a yw'n ymadrodd pert nad oes neb ond dysgwyr yn ei ddefnyddio?* Byddai corpws yn gallu rhoi canllaw o ran y penderfyniadau hyn ac eraill, yn hytrach na'n greddfau.

Rwy'n fy ystyried fy hun yn siaradwr Cymraeg iaith gyntaf, ar ôl cael fy magu yn Gymraeg mewn ardal Gymraeg yn ne orllewin Cymru. Fodd bynnag, mae ein cysyniadau am bwy neu beth yw iaith 'gyntaf' ac 'ail' iaith yn newid yn gyflym. Byddai cael darlun cynhwysfawr o sut y defnyddir yr iaith, yn arbennig yr iaith lafar, yn ddefnyddiol wrth ailddiffinio'r cysyniadau hyn, a sut byddwn ni'n mynd ati i ddysgu, addysgu ac asesu yn Gymraeg. Mae angen yr holl adnoddau a'r holl ymchwilwyr posibl er mwyn archwilio'r posibiladau.

*Dr Emyr Davies, [emyr.davies@cbac.co.uk](mailto:emyr.davies@cbac.co.uk)*

## CorCenCC ar-lein

Mae'r wybodaeth ddiweddaraf am ddatblygiadau'r prosiect hefyd ar gael drwy Facebook [www.facebook.com/CorCenCC/](http://www.facebook.com/CorCenCC/); Twitter <https://twitter.com/corcencc> (gallwch ein trydar @CorCenCC). Gallwch hefyd gysylltu â ni drwy anfon neges i gyfeiriad ebost y prosiect: [corcencc@caerdydd.ac.uk](mailto:corcencc@caerdydd.ac.uk) neu ewch i'n gwefan, sef: <http://sites.cardiff.ac.uk/corcencc/>

## CorCenCC mewn cynadleddau (Mai)

- **Piao, S., Rayson, P.,** Archer, D., Bianchi, F., Dayrell, C. El-Haj, M., Jiménez R-M., **Knight, D.,** Michal Křen, M., Löfberg L., Nawab, R., Shafi, J., The, P-L. a Mudraya, O. (2016). Lexical Coverage Evaluation of Large-scale Multilingual Semantic Lexicons for Twelve Languages. Papur i'w gyflwyno yng nghynhadledd *LREC (Language Resources Evaluation) 2016*, Mai 2016, Slofenia.

Mae CorCenCC yn brosiect ymchwil a ariennir gan ESRC/AHRC (Grant Rhif ES/M011348/1). Mae tîm CorCenCC yn cynnwys y **Prif Ymchwilydd**-Dawn Knight; y **Cyd-Ymchwilwyr** - Tess Fitzpatrick, Steve Morris, Irena Spasic, Paul Rayson, Enlli Thomas, Mark Stonelake a Jeremy Evas; y **Cynorthwywyr Ymchwil** - Steven Neale, Jennifer Needs, Mair Rees, Scott Piao a Gareth Watkins; **Ymgynghorwyr** - Kevin Donnelly, Kevin Scannell, Laurence Anthony, Tom Cobb, Michael McCarthy a Margaret Deuchar; **Grŵp Ymgynghorol y Prosiect** - Karen Corrigan, Llion Jones, Maggie Tallerman, Mair Parry-Jones a Gwen Awbery, Emyr Davies (CBAC-WJEC), Gareth Morlais (Llywodraeth Cymru), Owain Roberts (Llyfrgell Genedlaethol Cymru), Aran Jones (Saysomethingin.com) ac Andrew Hawke (Geiriadur Prifysgol Cymru).

Os oes gennych unrhyw sylwadau neu gwestiynau am gynnwys y cylchlythyr hwn, cysylltwch â Dr Dawn Knight: [KnightD5@caerdydd.ac.uk](mailto:KnightD5@caerdydd.ac.uk)



Arts & Humanities  
Research Council



### Greetings from the PI

*The third month of the CorCenCC project is now in full swing and it is great to report that work is going well. We are beginning to make some significant progress on individual work packages (WPs) on the project, and I am proud to say that we have already achieved a lot in this relatively short space of time. In this month's newsletter we will be providing brief updates on some of the on-going developments across specific WPs of the project and will be showcasing, in particular, recent work on the development of a Welsh semantic tagger that is being led by the team at Lancaster University (WP3). In addition to this, this month sees the inclusion of a new feature which will introduce you to individual members of the CorCenCC project team.*

*Happy reading, Dr Dawn Knight (Cardiff University)*

### Updates from CorCenCC Work Packages (WPs) 1, 2, 4, 5

Work on the CorCenCC project is distributed across 6 coordinated work packages, each with specific tasks, aims and objectives. WP0 involves on-going design, scoping and training activities, and involves all members of the project team. Brief updates on WPs 1, 2, 3, 4 are provided below (see the next section for updates on WP3).



#### **Aim: Collect, transcribe and anonymise the data (lead: Steve Morris)**

The WP1 team have been developing the sampling frame for data to be collected for CorCenCC. Following extensive research, a draft version of this has been drawn up and circulated to Welsh language and corpus experts for feedback. The team have also been working on devising transcription and anonymization conventions for the corpus.



#### **Aim: Develop the part-of-speech tag-set/tagger (lead: Dawn Knight)**

Progress continues on our bespoke part-of-speech tagset and tagging tools for Welsh. Additionally, we are putting plans in place for the production in the coming months of a gold-standard dataset for training/evaluating Welsh natural language processing tools!



#### **Aim: Scope/construct the online pedagogic toolkit (leads: Enlli Thomas/Tess Fitzpatrick)**

Work has begun on exploring the kinds of online tools already available to Welsh learners to avoid duplication of existing resources, and plans are underway to conduct a survey to investigate what resources learners and teachers already use and what they would ideally like to see developed as part of CorCenCC's pedagogical toolkit.



#### **Aim: Construct infrastructure to host CorCenCC and build the corpus (lead: Irena Spasic)**

Our project server is almost up and running and work is beginning on implementing our data crowdsourcing application, the first big steps in putting the CorCenCC infrastructure in place (keeping checking [www.corcenc.org](http://www.corcenc.org) for developments).



## Developing Semantic Annotation Tool for Welsh Language Analysis

In the CorCenCC Project, which aims to construct a large Welsh corpus and develop Welsh language tools, we are developing a suite of Welsh language processing software to assist the analysis and search for various information stored in the corpus data.



One of the major tools under development in this project at Lancaster University is the semantic annotation software, which aims to help to carry out the analysis of Welsh language data at the semantic level in a large scale. This tool will process Welsh language samples automatically and label each word or phrase with its coarse-grained meaning. Based on and extending the Lancaster USAS system (<http://ucrel.lancs.ac.uk/usas/>), it will be capable of assigning semantic categories, in the form of tags such as I1.3 (Arian: Prisiau/Money: Price), K5.1 (Chwaraeon/Sports), to Welsh words, idioms and other fixed phrases (also called multi-word expressions) based on a semantic classification scheme. This scheme consists of 232 semantic categories falling under 21 major categories, as shown below:

Tag	Definition	Tag	Definition
A	TERMAU CYFFREDINOL A HANIAETHOL (GENERAL AND ABSTRACT TERMS)	N	RHIFAU A MESUR (NUMBERS AND MEASUREMENT )
B	B Y CORFF A'R UNIGOLYN (THE BODY & THE INDIVIDUAL)	O	SYLWEDDAU, DEFNYDDIAU, GWRTHRYCHAU AC OFFER (SUBSTANCES, MATERIALS, OBJECTS AND EQUIPMENT)
C	CELFA CHREFFT (ARTS AND CRAFTS)	P	ADDYSG (EDUCATION)
E	GWEITHREDIADAU, CYFLYRAU A PHROSESAU EMOSIYNOL (EMOTIONAL ACTIONS, STATES & PROCESSES)	Q	GWEITHREDIADAU, CYFLYRAU A PHROSESAU IEITHYDDOL (LINGUISTIC ACTIONS, STATES AND PROCESSES )
F	BWYD A FFERMIO (FOOD & FARMING)	S	GWEITHREDIADAU, CYFLYRAU A PHROSESAU CYMDEITHASOL (SOCIAL ACTIONS, STATES AND PROCESSES)
G	LLYWODRAETH A'R CYHOEDD (GOVERNMENT AND THE PUBLIC DOMAI)	T	AMSER (TIME )
H	PENSAERNIAETH, ADEILADAU, TAI A'R CARTREF (ARCHITECTURE, BUILDINGS, HOUSES & THE HOME)	W	Y BYD A'N HAMGYLCHEDD (THE WORLD AND OUR ENVIRONMENT)
I	ARIAN A MASNACH (MONEY & COMMERCE)	X	GWEITHREDIADAU, CYFLYRAU A PHROSESAU SEICOLEGOL (PSYCHOLOGICAL ACTIONS, STATES AND PROCESSES)
K	ADLONIAN, CHWARAEON A GEMAU (ENTERTAINMENT, SPORTS AND GAMES)	Y	GWYDDONIAETH A THECHNOLEG (SCIENCE AND TECHNOLOGY)
L	BYWYD A PHETHAU BYW (LIFE AND LIVING THINGS )	Z	ENWAU A GEIRIAU GRAMADEGOL (NAMES AND GRAMMATICAL WORDS)
M	SYMUD, LLEOLIAD, TEITHIO A CHLUDIANT (MOVEMENT, LOCATION, TRAVEL AND TRANSPORT)		

It is a highly challenging task to develop an accurate semantic annotation tool. Our tool will combine a Welsh semantic lexical knowledge base (in essence a large machine readable dictionary) and a range of **word sense disambiguation** methods to achieve a high accuracy of the semantic annotation. Furthermore, we will explore crowdsourcing approaches, in which Welsh speaking communities will be invited to participate, to help to construct the semantic lexical knowledge base on a large scale and to improve the performance of the tool.

The Welsh semantic annotation tool can be useful for a variety of academic and practical purposes. Firstly, the entire Welsh corpus to be constructed in the CorCenCC Project will be tagged using this tool, which will facilitate various semantic analysis of the corpus data and a reliable and fast search of the corpus data by various semantic categories, as demonstrated by the Lancaster Wmatrix corpus research site (<http://ucrel.lancs.ac.uk/wmatrix/>). In addition, the semantic annotation tool can be used to improve Welsh language based ICT systems. For example, it can allow an ICT system to rapidly extract and analyse specific semantic information from online Welsh language sources on a large scale and help to link useful online sources for users.



Although the CorCenCC Project started only recently, the development of the semantic annotation tool is progressing well and we have released the translated set of categories as shown above.

***Dr Paul Rayson (WP3 lead) and Dr Scott Piao (RA), Lancaster University***

## **Meet the team**

*Every month we will be featuring a different member of the extended CorCenCC team in our newsletter. This will give everyone a chance to tell you a little about their background; what they want to see from CorCenCC and how they think it might contribute to their own work, or more broadly, to the work of others in Wales. This month we spotlight Dr Emyr Davies from CBAC-WJEC, who is part of the CorCenCC Project Advisory Group.*

### **Profile: Dr Emyr Davies**

At last! That was my reaction to the news that there would be a comprehensive corpus of the Welsh language. There have been small corpora developed in the past, but nothing on this scale with so many potential uses. I've been interested in language and linguistics for a long time. After graduating in Welsh from the University of Wales Aberystwyth and gaining my PhD, I spent 11 years teaching Welsh to adults in Trinity College, Carmarthen. I spent a lot of time developing new resources, but this was before the digital revolution; ideas such as electronic corpora simply didn't exist. I took a year off in 1996 to study for a Higher Diploma in General Linguistics at UCD Dublin, which I enjoyed very much. It's a useful grounding for any one working in languages, even if they then go on to work in some aspect of applied linguistics, which is exactly what I did.



I left Carmarthen to join CBAC-WJEC, the exam board based in Cardiff, in 2001 and have been working primarily in assessment and developing resources for adult learners of Welsh since then. The main part of my job has been to establish a suite of language exams for the sector, which are now in place. Another part of my job has been to represent the Welsh language exams in ALTE, the Association of Language Testers in Europe.

This has been a great experience for me and given me opportunities to attend workshops and learn from some of the leading experts in the field. Not many people work purely on assessment, certainly in Welsh, and ALTE (for one thing) provides a professional community of organisations focused on language testing. ALTE also requires its members to be audited, or to provide evidence that we reach a high quality profile. The same

robust standards apply to all languages and organisations who wish to be full members, including all major European languages and many less widely used languages, including Basque, Irish and Welsh.

ALTE's approach is based on an evidence-based approach. In other words, can we provide evidence to the support the decisions we make? A comprehensive corpus could be useful in developing tests, as well as resources for learning, and would provide us with evidence to make a thousand and more decisions. Many questions spring to mind: *How do we choose aspects of language to teach (and test) first? Do frequency tables for Welsh vocabulary correspond to English? Are there collocations which we should include in our teaching materials? How can we include 'redundant' language, e.g. Well... uhm... in an authentic way? Which aspects of dialect are most useful for learners? Is this or that phrase actually useful, or is it simply a quaint expression that nobody uses any more except learners?* A corpus could guide these and other decisions, rather than our intuitions.

I consider myself to be a first language Welsh speaker, having been raised in Welsh in a Welsh speaking area of south west Wales. However, our notions of who or what is a 'first' language and 'second' language user are changing fast. Having a comprehensive picture of how the language is used, especially spoken language, will be useful in redefining these notions, and how we approach learning, teaching and assessment in Welsh. We need all the resources and researchers we can find in order to explore the possibilities.

*Dr Emyr Davies, [emyr.davies@cbac.co.uk](mailto:emyr.davies@cbac.co.uk)*

## CorCenCC online

You can keep up to date with developments on the project via Facebook [www.facebook.com/CorCenCC/](https://www.facebook.com/CorCenCC/); Twitter <https://twitter.com/corcencc> (Tweet us @CorCenCC). You can also contact us on the project email address: [corcencc@cardiff.ac.uk](mailto:corcencc@cardiff.ac.uk) or visit our holding website at: <http://sites.cardiff.ac.uk/corcencc/>

## CorCenCC @ conferences and events (May)

- **Piao, S., Rayson, P.,** Archer, D., Bianchi, F., Dayrell, C. El-Haj, M., Jiménez R-M., **Knight, D.,** Michal Křen, M., Löfberg L., Nawab, R., Shafi, J., The, P-L. and Mudraya, O. (2016). Lexical Coverage Evaluation of Large-scale Multilingual Semantic Lexicons for Twelve Languages. Poster presentation to be delivered at the *LREC (Language Resources Evaluation) 2016 Conference*, May 2016, Slovenia.

CorCenCC is an ESRC/AHRC funded research project (Grant Number ES/M011348/1). The CorCenCC team includes **PI** - Dawn Knight; **CIs** - Tess Fitzpatrick, Steve Morris, Irena Spasic, Paul Rayson, Enlli Thomas, Mark Stonelake and Jeremy Evas; **RAs** - Steven Neale, Jennifer Needs, Mair Rees, Scott Piao and Gareth Watkins; **Consultants** - Kevin Donnelly, Kevin Scannell, Laurence Anthony, Tom Cobb, Michael McCarthy and Margaret Deuchar; **Project Advisory Group** - Karen Corrigan, Llion Jones, Maggie Tallerman, Mair Parry-Jones and Gwen Awbery, Emyr Davies (CBAC-WJEC), Gareth Morlais (Welsh Government), Owain Roberts (National Library of Wales), Aran Jones (Saysomethingin.com) and Andrew Hawke (University of Wales Dictionary of the Welsh Language).

If you have any comments or questions about the content of this newsletter please contact Dr Dawn Knight: [KnightD5@cardiff.ac.uk](mailto:KnightD5@cardiff.ac.uk)



Arts & Humanities Research Council