



**Universidade de Aveiro**  
2012

Secção Autónoma de Ciências da Saúde

**MARIA INÊS  
MORGADO  
OLIVEIRA  
MARTINS**

**EVOLUÇÃO E RELEVÂNCIA FUNCIONAL DOS  
PARÁLOGOS DO GENE *ATAXIN-3***

**EVOLUTION AND FUNCTIONAL RELEVANCE OF  
*ATAXIN-3* PARALOGUES**





**Universidade de Aveiro**  
2011

Secção Autónoma de Ciências da Saúde

**MARIA INÊS  
MORGADO  
OLIVEIRA  
MARTINS**

**EVOLUÇÃO E RELEVÂNCIA FUNCIONAL DOS  
PARÁLOGOS DO GENE *ATAXIN-3***

**EVOLUTION AND FUNCTIONAL RELEVANCE OF  
*ATAXIN-3* PARALOGUES**

Dissertação apresentada à Universidade de Aveiro para cumprimento dos requisitos necessários à obtenção do grau de Mestre em Biomedicina Molecular, realizada sob a orientação científica da Doutora Sandra Cristina da Silva Martins, Investigadora do Instituto de Patologia e Imunologia Molecular da Universidade do Porto, e a co-orientação da Doutora Ana Gabriela da Silva Cavaleiro Henriques, Professora Auxiliar Convidada da Secção Autónoma das Ciências da Saúde da Universidade de Aveiro.

Esta dissertação foi desenvolvida no Instituto de Patologia e Imunologia Molecular da Universidade do Porto (IPATIMUP), um laboratório associado do Ministério Português da Educação e Ciência.



## **o júri**

Presidente

**Prof. Doutora Odete Abreu Beirão da Cruz e Silva**  
Professora Auxiliar com Agregação da Secção Autónoma das Ciências da Saúde da Universidade de Aveiro

**Doutora Sandra Cristina da Silva Martins**  
Investigadora no Instituto de Patologia e Imunologia Molecular da Universidade do Porto

**Prof. Doutora Ana Gabriela da Silva Cavaleiro Henriques**  
Professora Auxiliar Convidada da Secção Autónoma das Ciências da Saúde da Universidade de Aveiro

**Doutora Ana Goios Borges de Almeida**  
Investigadora no Instituto de Patologia e Imunologia Molecular da Universidade do Porto



## **agradecimentos**

Ao Professor António Amorim pela oportunidade de integração no seu entusiasmante grupo de investigação.

À Doutora Sandra Martins por me ter acompanhado incansavelmente, pela sua disponibilidade, paciência e compreensão. Agradeço muito também por me ter inspirado e mostrado toda a sua dedicação e motivação.

À Professora Ana Gabriela Henriques, por desde sempre me ter acompanhado e mostrado os caminhos certos a seguir.

Ao IPATIMUP pelo financiamento, pelo acolhimento e apoio científico sem o qual o sucesso do trabalho experimental não teria sido possível.

A todos os elementos do Grupo de Genética Populacional, que tão bem me acolheram.

À Doutora Odete da Cruz e Silva por me ter proporcionado estes excelentes 5 anos de aprendizagem, bem como todas as oportunidades únicas que tive e me fizeram crescer, quer a nível profissional como pessoal. Agradeço ainda a todos os Professores/Tutores que me acompanharam e ajudaram da melhor maneira possível (principalmente às professoras Sandra Vieira, Ana Gabriela Henriques, Margarida Fardilha e Sandra Rebelo).

À Catarina Seabra, pelo constante e incondicional apoio, conselhos e presença nos momentos bons e menos bons.

Ao Semião por, apesar da distância, ser um grande amigo e estar sempre lá, pronto a ouvir-me nos momentos de felicidade e tristeza.

À Joana Tavares, Chita, Ana Sol e Cátia pela amizade, carinho e apoio. Obrigada ainda por todos os ótimos momentos partilhados, os quais nunca esquecerei.

À Catarina Xavier, Filipa, Marisa, Sofia e Alexandre pelos ótimos momentos passados e o apoio constante ao longo de todo o ano.

A todos os meus colegas de Ciências Biomédicas de Aveiro, que me fizeram crescer e aprender a cada passo.

Aos meus amigos de Aveiro, principalmente à Raleira, ao Recha, Tubrinho, Badaró, Big, Renato, Baka, Jorge e ao David, que me acompanharam desde cedo, e me deram a amizade e os melhores anos de Universidade que mais poderia desejar!

Ao João Fernandes por sempre me apoiar e animar nos momentos mais difíceis, e principalmente, por ser tão importante e me dar tantos dos melhores momentos da minha vida.

À minha irmã Mariana, por ter estado sempre ao meu lado, apoiado e ouvido, principalmente nos momentos em que mais precisava que me ouvissem.

Aos meus pais, pelo amor, ensinamentos sábios e pelo incansável apoio e motivação que sempre me deram ao longo da vida.

Às minhas avós por todo o amor, orgulho e interminável apoio.





## palavras-chave

Ataxina-3, Retrotransposição, Genes parálogos, Evolução, Doença de Machado-Joseph

## resumo

O gene *ataxin-3* (*ATXN3*; 14q32.1) codifica uma proteína expressa ubiquamente, envolvida na via ubiquitina-proteassoma e na repressão da transcrição. Grande relevância tem sido dada ao gene *ATXN3* após a identificação de uma expansão  $(CAG)_n$  na sua região codificante, responsável pela ataxia mais comum em todo o mundo, SCA3 ou doença de Machado-Joseph (DMJ). A DMJ é uma doença neurodegenerativa, autossômica dominante, de início tardio. O tamanho do alelo expandido explica apenas uma parte do pleomorfismo da doença, evidenciando a importância do estudo de outros modificadores. Em doenças de poliglutaminas (poliQ), a toxicidade é causada por um ganho de função da proteína expandida; no entanto, a proteína normal parece ser, também, um dos agentes modificadores da patogênese. O gene *ATXN3* possui dois parálogos humanos gerados por retrotransposição: *ataxin-3 like* (*ATXN3L*) no cromossoma X, e *LOC100132280*, ainda não caracterizado, no cromossoma 8. Estudos *in vitro* evidenciaram a capacidade da *ATXN3L* para clivar cadeias de ubiquitina, sendo o seu domínio proteolítico mais eficiente do que o domínio da *ATXN3* parental.

O objetivo deste estudo foi explorar a origem e a evolução das retrocópias *ATXN3L* e *LOC100132280* (aqui denominadas *ATXN3L1* e *ATXN3L2*), assim como testar a relevância funcional de ambas através de abordagens evolutivas e funcionais. Deste modo, para estudar a divergência evolutiva dos parálogos do gene *ATXN3*: 1) analisaram-se as suas filogenias e estimou-se a data de origem dos eventos de retrotransposição; 2) avaliaram-se as pressões seletivas a que têm sido sujeitos os três parálogos, ao longo da evolução dos primatas; e 3) explorou-se a evolução das repetições CAG, localizadas em três contextos genômicos diferentes, provavelmente sujeitos a diferentes pressões seletivas. Finalmente, para o retrogene que conserva uma *open reading frame* (ORF) intacta, *ATXN3L1*, analisou-se, *in silico*, a conservação dos locais e domínios proteicos da putativa proteína. Ademais, para este retrogene, foi estudado o padrão de expressão de mRNA, através da realização de PCR de Transcriptase Reversa, em 16 tecidos humanos.

Os resultados obtidos sugerem que dois eventos independentes de retrotransposição estiveram na origem dos retrogenes *ATXN3L1* e *ATXN3L2*, tendo o primeiro ocorrido há cerca de 63 milhões de anos (Ma) e o segundo após a divisão Placental-Macarrínios, há cerca de 35 Ma. Adicionalmente, outras retrocópias foram encontradas em primatas e outros mamíferos, correspondendo, no entanto, a eventos mais recentes e independentes de retrotransposição. A abordagem evolutiva mostrou a existência de algumas restrições seletivas associadas à evolução do gene *ATXN3L1*, à semelhança do que acontece com *ATXN3*. Por outro lado, *ATXN3L2* adquiriu códons *stop* prematuros que, muito provavelmente, o tornaram num pseudogene processado. Os resultados da análise de expressão mostraram que o gene *ATXN3L1* é transcrito, pelo menos, em testículo humano; no entanto, a otimização final da amplificação específica dos transcriptos *ATXN3L1* permitirá confirmar se a expressão se estende a outros tecidos. Relativamente ao mecanismo de mutação inerente à repetição CAG, os dois parálogos mostraram diferentes padrões de evolução: a retrocópia *ATXN3L1* é altamente interrompida e pouco polimórfica, enquanto a *ATXN3L2* apresenta tratos puros de  $(CAG)_n$  em algumas espécies e tratos hexanucleotídicos de CGGCAG no homem e no chimpanzé. A recente aquisição da repetição CGGCAG pode ter resultado de uma mutação inicial

de CAG para CGG, seguida de instabilidade que proporcionou a expansão dos hexanucleótidos. Estudos futuros poderão ser realizados no sentido de confirmar o padrão de expressão do gene *ATXN3L1* e de detetar proteína endógena *in vivo*. Adicionalmente, a caracterização da proteína ataxina-3 like 1 e dos seus interatores moleculares poderá providenciar informação acerca da sua relevância no estado normal e patológico.

## keywords

Ataxin-3, Retrotransposition, Paralogue genes, Evolution, Machado-Joseph disease.

## abstract

*Ataxin-3* gene (*ATXN3*; 14q32.1) encodes a ubiquitously expressed protein involved in the ubiquitin-proteasome pathway and in transcription repression. Much attention has been given to *ATXN3* since the identification of an expanded (CAG)<sub>n</sub> tract in its coding region, responsible for the most common dominant ataxia worldwide, SCA3 or Machado-Joseph disease (MJD). MJD is an autosomal dominant, late-onset neurodegenerative disorder. The size of the expanded allele explains only part of the disease pleomorphism, highlighting the importance of studying other modifiers. In polyglutamine (polyQ) diseases, toxicity is caused by gain of function of the expanded protein; however, the normal protein appears to be one of the pathogenesis-modifying agents.

The gene *ataxin-3* has two human paralogues, generated by retrotransposition: *ataxin-3 like* (*ATXN3L*) on the X chromosome of humans and the yet uncharacterized *LOC100132280* on chromosome 8. An *in vitro* study showed the ability of *ATXN3L* to cleave ubiquitin from their substrates, with its proteolytic domain even more efficient than the domain of the parental *ATXN3*.

The aim of this study was to explore the origin and evolution of *ATXN3L* and *LOC100132280* retrocopies (here named *ATXN3L1* and *ATXN3L2*) and to test the functional relevance of both human paralogues through evolutionary and functional approaches. Thus, to study the evolutionary divergence of *ATXN3* paralogues, we have 1) analysed their phylogeny and estimated the time of the retrotransposition events; 2) assessed the selective constraints that have been underlying all the three paralogues over primate evolution; and 3) explored the evolution of the (CAG)<sub>n</sub> tracts placed on three different chromosomal backgrounds within paralogue genes that have, probably, been under different selective pressures. Finally, for the retrogene that conserved the ORF, *ATXN3L1*, we analyzed the protein domain conservation and the mRNA expression pattern by performing Reverse Transcriptase-PCR in 16 human tissues.

Our results suggested that two independent retrotransposition events have been on the origin of *ATXN3L1* and *ATXN3L2*, the first occurred about 63 million years ago (MYA) and the second after the Platyrrhini-Catarrhini split, about 35 MYA. In addition, several other retrocopies have been found in primates and other mammals, but additional independent and younger retrotransposition events seemed to be on their origin. Our evolutionary studies suggested that *ATXN3L1* has been under some selective constraints over primate evolution, as the parental *ATXN3*. *ATXN3L2* gained premature stop codons that seem to have turned it into a pseudogene. In addition, we confirmed that *ATXN3L1* is a transcriptionally-active retrogene since we observed mRNA expression, at least, in human testis. A refined optimization of the methodology to specifically amplify *ATXN3L1* cDNA will be, however, necessary to assess the complete expression pattern of this retrogene. As for the (CAG)<sub>n</sub> tract, the *ATXN3* paralogues presented different evolutionary patterns: *ATXN3L1* showed a poorly polymorphic and highly interrupted tract across the primate lineage, whereas *ATXN3L2* presented a pure (CAG)<sub>n</sub> in some species and a polymorphic hexanucleotide repeat in humans and chimpanzees. This recent acquisition of a repetitive CGGCAG resulted, probably, from a CAG to CCG mutation followed by instability that encompassed the six bases instead of the CAG alone.

Future studies may be performed in order to confirm the expression pattern of *ATXN3L1* and to detect the endogenous protein *in vivo*. Further characterization of *ataxin-3 like 1* and its molecular interactors will give us insight into the actual relevance of this protein in normal and/or pathological states.



# Contents

---

<b>Abbreviations .....</b>	<b>3</b>
<b>Chapter 1 - Introduction .....</b>	<b>5</b>
<b>1.1 - Mechanisms of retroposition .....</b>	<b>8</b>
<b>1.2- Machado-Joseph disease .....</b>	<b>10</b>
1.2.1 - Clinical presentation and epidemiology .....	10
1.2.2 - Molecular genetics and pathogenesis .....	11
<i>ATXN3</i> gene .....	11
Mechanisms of <i>ATXN3</i> mutation .....	12
Ataxin-3 protein and its physiologic role .....	13
Polyglutamine expansion of <i>ATXN3</i> and neural cell death .....	14
<b>1. 3 – Studies on the relevance of <i>ataxin</i> paralogues in SCAs .....</b>	<b>15</b>
<b>Aims .....</b>	<b>17</b>
<b>Chapter 2 - Material and methods .....</b>	<b>19</b>
Subjects .....	21
<b>Part 1 – Evolutionary approach .....</b>	<b>22</b>
Compilation and alignment of <i>ATXN3L1</i> and <i>ATXN3L2</i> sequences .....	22
Phylogenetic trees and genetic distances .....	23
Synteny of <i>ATXN3</i> , <i>ATXN3L1</i> and <i>ATXN3L2</i> .....	23
dN/dS ratio (omega) calculations .....	23
Primer design, DNA amplification and sequencing .....	24
Nucleotide diversity of human <i>ATXN3</i> retrocopies .....	28
(CAG) <sub>n</sub> tract analysis in <i>ATXN3</i> , <i>ATXN3L1</i> and <i>ATXN3L2</i> .....	28
<b>Part 2 – Functional approach .....</b>	<b>28</b>
ORF prediction for <i>ATXN3L1</i> and <i>ATXN3L2</i> .....	28

Conservation of the ATXN3L1 putative protein .....	28
mRNA expression of <i>ATXN3L1</i> .....	29
<b>Chapter 3 - Results .....</b>	<b>32</b>
<b>Part 1 – Evolutionary history of <i>ATXN3</i> paralogues .....</b>	<b>33</b>
1. Onset of <i>ATXN3</i> retrocopies.....	33
1.1. Identification of <i>ATXN3L1</i> and <i>ATXN3L2</i> orthologues in the primate lineage... 33	
1.2. Phylogeny of <i>ATXN3</i> paralogues in primates .....	39
1.3. Identification of <i>ATXN3</i> retrocopies in other mammals .....	40
2. <i>ATXN3</i> transcripts involved in the retrotransposition events.....	45
3. Selective signatures underlying <i>ATXN3L1</i> and <i>ATXN3L2</i> .....	45
4. Nucleotide diversity of <i>ATXN3L1</i> and <i>ATXN3L2</i> .....	47
5. Evolution of the (CAG) <sub>n</sub> tract in <i>ATXN3</i> , <i>ATXN3L1</i> and <i>ATXN3L2</i> .....	48
<b>Part 2 – Insights into <i>ATXN3L1</i> and <i>ATXN3L2</i> functional relevance .....</b>	<b>50</b>
1. ORF predictions for <i>ATXN3L1</i> and <i>ATXN3L2</i> .....	50
2. Analysis of <i>ATXN3L1</i> putative protein domains in comparison the parental <i>ATXN3</i> 52	
3. mRNA expression profile of <i>ATXN3L1</i> in humans .....	53
<b>Chapter 4 - Discussion .....</b>	<b>55</b>
<b>Chapter 5 - Final remarks and future perspectives.....</b>	<b>63</b>
<b>References.....</b>	<b>69</b>
<b>Appendix .....</b>	<b>73</b>

# Abbreviations

---

- AO** – Age-at-onset
- ATXN1** – Ataxin-1
- ATXN1L** – Ataxin-1 like
- ATXN3** – Ataxin-3
- ATXN3L1** – Ataxin-3 like 1
- ATXN3L2** – Ataxin-3 like 2
- BLAST** - Basic local alignment search tool
- BLAT** - BLAST-like alignment tool
- bp** – base pairs
- CIC** – Capicua
- DUB** – Deubiquitinating enzyme
- JD** – Josephin domain
- JOSD1** – Josephin 1
- JOSD2** – Josephin 2
- L1** – Long interspersed nuclear element 1
- LINE** – Long interspersed nuclear element
- LTR** – Long terminal repeats
- MJD** – Machado-Joseph disease
- MYA** – Million years ago
- NCBI** – National Center for Biotechnology Information
- NES** – Nuclear export signal
- NI** – Nuclear inclusions
- NLS** – Nuclear localization signal
- ORF** – Open reading frame
- PolyQ** – poly glutamine
- RT-PCR** – Reverse transcriptase polymerase chain reaction
- SCA1** – Spinocerebellar ataxia type 1
- SCA3** – Spinocerebellar ataxia type 3
- SNP** – Single nucleotide polymorphism

University of Aveiro – Master's in Molecular Biomedicine  
IPATIMUP  
2012

**STR** – Simple tandem repeat

**SVA** – (SINE/VNTR/*Alu*)

**Ub** – Ubiquitin

**UCSC** –University of California, Santa Cruz

**UIM** – Ubiquitin interacting motif

**UTR** – Untranslated region



## Chapter 1

# Introduction

---



## Chapter 1

# Introduction

---

New genes are thought to contribute to the origin of adaptive evolutionary innovations and, thus, to lineage- or species-specific phenotypic traits. Traditionally, the origin of these new genes is associated with gene duplication, but other mechanisms have recently received increasing attention, such as retrotransposition, which may play an important role in mammal genome evolution.<sup>1, 2, 3</sup> Retrotransposition consists in the re-integration of reverse transcribed mRNA molecules in the genome.<sup>4</sup> At the sequence level, retroposed copies (retrocopies) can be intact or not, depending on the presence of mutations or premature stop codons. Still, if the open reading frame is conserved, retrocopies lack many of the genetic features of their parental genes (such as introns and regulatory elements) and, for this reason, they have been considered as confounding factors for a long time. Most retrocopies turned into pseudogenes in mammals, but some of them have recruited upstream regulatory elements and became functional (these are commonly called retrogenes).<sup>5</sup> Indeed, many recent studies have shown that a larger part of retrotransposons than the previously predicted are functional, and that they can modulate other pre-existing cellular factors.<sup>1,5</sup>

This thesis project is focused on the study of two retrocopies of *ataxin-3* (*ATXN3*), a ubiquitously expressed gene that codes a protein mainly involved in deubiquitination and in transcription repression. *ATXN3* is also the gene responsible for Machado-Joseph disease (MJD), when its coding repetitive CAG tract is expanded above 61 repeats. MJD or spinocerebellar ataxia type 3 (SCA3) is a late-onset neurodegenerative disorder, characterized by a large pleomorphism. Thus, the search for disease modifiers has been a major point of interest in this disorder.

The origin, features and functional relevance of the poorly studied *ATXN3* retrocopies, *ATXN3L* and *LOC100132280* (here named *ATXN3L1* and *ATXN3L2*, respectively) remain completely unexplored until these days. With this project we aimed at characterizing *ataxin-3* paralogues in the human genome and throughout the evolution to gain insight into the potential roles that they have been playing.

## 1.1 - Mechanisms of retroposition

The key retrotransposition enzyme stems from different types of retrotransposable elements, depending on the organism. These elements can be subdivided in two groups, distinguished by the presence or absence of long terminal repeats (LTRs). In humans, LTR elements are endogenous retroviruses (HERVs), presently with very limited activity, which accounts for only ~8% of the genome. On the other hand, human non-LTRs include long interspersed nuclear element 1 (LINE 1 or L1), *Alu* and SVA (an element composed of a short interspersed region, a variable number of tandem repeats region and an *Alu*-like region - SINE/VNTR/*Alu*) elements, which collectively account for approximately one-third of the human genome (Figure 1).<sup>6</sup>

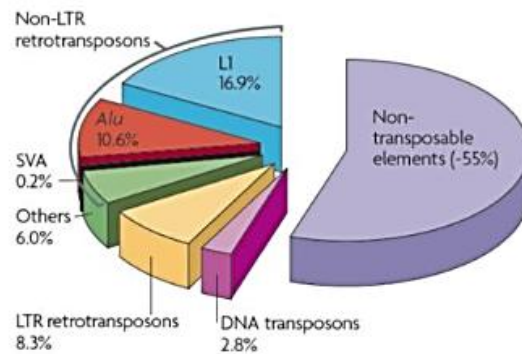
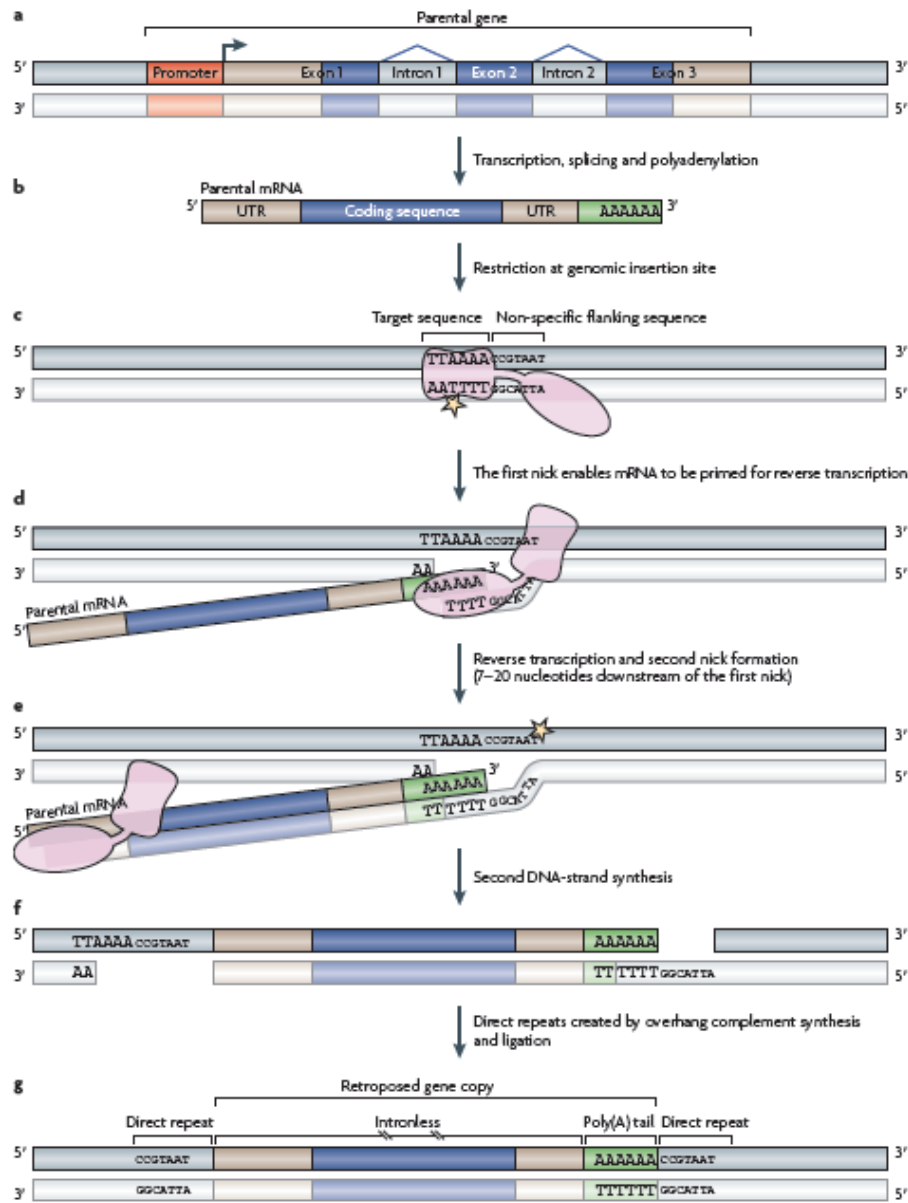


Figure 1 - Transposable and non-transposable element content of the human genome. About 33.7% of the human genome can currently be recognized as being derived from non-LTR retrotransposable elements (Cordaux et al., 2009).<sup>6</sup>

L1, widely present in mammals (about 25% of their genome), are responsible for a burst on the number of retrocopies in mammalian lineages.<sup>2</sup> These retrotransposable elements possess a reverse transcriptase and a endonuclease activity that can recognize any polyadenylated mRNA.<sup>7</sup> After gene transcription, the respective RNA is reverse-transcribed into DNA by the L1 elements, and the resulting cDNA is inserted into the genome at a new location (Figure 2).<sup>4,7</sup>



**Figure 2 - Retrotransposition mechanism:** a) retroposition is initiated with the transcription of a parental gene by RNA polymerase; b) a mature mRNA is produced, with the processing of the resulting RNA (by splicing and polyadenylation); c) L1 endonuclease domain (pink rectangle) mediates retroposition, by creating a first nick (yellow star) at the genomic site of insertion at the TTA AAA target sequence; d) this nick enables the mRNA to be primed for reverse transcription by the L1 reverse transcriptase domain (pink oval), which uses the parental mRNA as a template; e) second-strand nick generation, f) second DNA-strand synthesis; g) cDNA synthesis in the overhang regions created by the two nicks. (adapted from Kaessmann, 2009)<sup>1</sup>

To become expressed at a significant level and in a meaningful way, a new retrogene needs to obtain a core promoter and probably other elements, such as enhancers, that regulate its expression.<sup>1</sup> Generally, the expression of a retrocopy might benefit from: 1) its insertion into intronic sequences of host genes, enabling it to be integrated into new splice variants of the host gene; 2) its insertion into actively transcribed regions with an open chromatin structure, as this increases accessibility for the transcriptional machinery; 3) the recruitment of distant promoters in the genomic neighbourhood via the acquisition of a new untranslated exon–intron structure; 4) the recruitment of proto-promoters from retrotransposons or CpG island; 5) the inheritance of parental promoters through alternative transcriptional start sites used by the parental gene; and 6) *de novo* promoter evolution in the upstream flanking region of the insertion site by single nucleotide substitutions.<sup>1, 5</sup> Since retroposed copies often need to recruit regulatory elements to become transcribed, they are prone to accumulate genetic variants and to evolve new expression patterns relatively to the traditional gene duplication. As a consequence, novel functional roles can emerge from the transcription of new formed retrogenes.<sup>1</sup> Studies have revealed that retrogenes seem to evolve biased functions related to male functions, but others functions have also been described, for example, related to the brain.<sup>8, 9</sup>

In addition, to be heritable and hence of evolutionary relevance, retrotransposition needs to occur in the germ line. The fact that retrotransposition relies on duplication through an mRNA intermediate implies that only genes expressed in the germ line can be the source of new retrocopies.<sup>1</sup>

## **1.2- Machado-Joseph disease**

### **1.2.1 - Clinical presentation and epidemiology**

Machado-Joseph disease (MJD, also called spinocerebellar ataxia type 3, SCA3) is the most common type of spinocerebellar ataxia. SCAs are a large and complex group of late-onset diseases, characterized by progressive cerebellar dysfunction of afferent and efferent pathways, variably associated with other symptoms of the central and peripheral nervous systems. Other nervous system structures affected include the basal ganglia, brainstem nuclei, pyramidal tracts and post superior column and anterior horn of the spinal cord, as well as peripheral nerves.<sup>10, 11</sup> Although there are sporadic forms of ataxia, the term SCA is most often used to refer the

hereditary forms, and in particular the autosomal dominant ataxias.<sup>8</sup> Nearly 30 subtypes of SCAs have been described, differently classified based on the underlying causative mutation. MJD is one of the polyglutamine diseases since the causative mutation is a CAG repeat expansion in *ATXN3* gene, which encodes a stretch of glutamine amino acids in the corresponding protein.<sup>12, 13</sup>

In the particular case of MJD, the major signs are the progressive cerebellar ataxia and pyramidal signs. Minor, but more specific, features are external progressive ophtalmoplegia, dystonia, intention fasciculation-like movements of facial and lingual muscles, as well as bulging eyes but these symptoms can vary among patients. The mean age at onset (AO) of the disease is around 40 years, although it can vary greatly, with extremes of 4 and 70 years.<sup>11</sup> The mean survival time is 21 years (ranging from 7 to 29).<sup>14</sup> Thus, MJD is characterized by a high degree of pleomorphism, not only in the variability of the AO, but also in the neurological signs presented and in the resulting degree of incapacity.<sup>10, 15</sup>

SCAs are considered rare disorders, with estimated prevalence described as varying from 0.3 to 2.0 per 100,000 individuals. The relative frequency of MJD among SCAs varies also largely among populations. It is higher in Brazil (69-92%),<sup>16</sup> Portugal (58-74%),<sup>17</sup> Singapore (53%),<sup>18</sup> China (48-49%),<sup>19</sup> the Netherlands (44%),<sup>20</sup> Germany (42%)<sup>21</sup> and Japan (28-63%)<sup>22</sup> but even within each country, the geographic distribution pattern of MJD is not homogeneous.

Haplotype-based studies have suggested that two independent-origin mutations may explain the current worldwide geographic distribution of MJD. The first occurred about 6000 years ago in Asia (TTACAC or Joseph lineage, observed in the 5 continents); and the most recent, with less than 2000 years old, is responsible for the presence of MJD in Portugal (together with the most ancient mutation) and in a few other populations mostly linked to Portugal (GTGGCA or Machado lineage).<sup>23</sup>

### **1.2.2 - Molecular genetics and pathogenesis**

#### ***ATXN3* gene**

*ATXN3* was located in the long arm of chromosome 14 (14q32.1), by Takiyama in 1993.<sup>24</sup> The genomic structure of *ATXN3* spans about 48 kb, with the most frequent transcript containing 11 exons.<sup>25</sup> Two additional exons, 6a and 9a, have recently been described (Figure 3).<sup>26</sup>

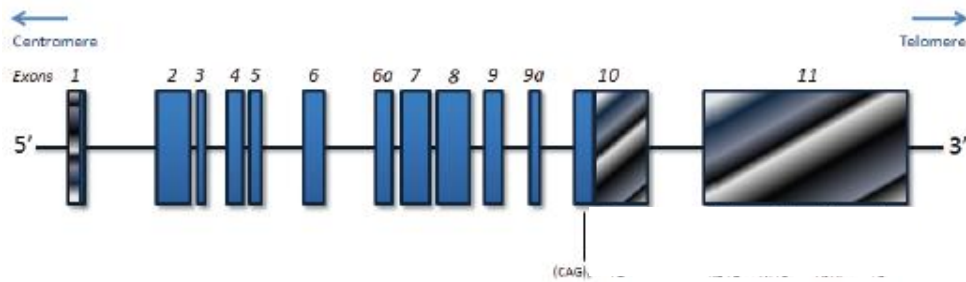


Figure 3 - *ATXN3* gene structure, representing the first described exons and the two new exons: 6a and 9a; the (CAG)<sub>n</sub> tract is localized in the exon 10. Transcribed regions are represented in blue, and not-transcribed and UTR regions in grey. (adapted from Bettencourt, 2011)<sup>13</sup>

The expansion of a (CAG)<sub>n</sub> tract within *ATXN3* exon 10 is the causative mutation of MJD.<sup>27</sup> Usually, wild-type alleles range from 12 to 44 CAG repeats, whereas expanded alleles comprise 61 to 87 CAGs. Intermediate size alleles are rare, and their role in disease presentation is still under debate.<sup>28, 29</sup> The size of the expansion has been found, until date, the parameter that best correlates with the AO of patients; however, it explains only 45-76% of AO variability suggesting a major influence of other factors.<sup>30, 31</sup> Gender was initially proposed to be one of these factors<sup>32</sup>, but this was later contradicted, as familial factors were suggested to mask the gender effect.<sup>33</sup> It was demonstrated, instead, that variation in AO accounted for an effect common to a sibship, independently of the CAG repeat length. Less controversial is the dosage effect on disease presentation: homozygous patients with two mutant alleles show, most frequently, a more severe disease phenotype, with earlier onset than those presenting only one mutated allele.<sup>34</sup>

The process underlying the (CAG)<sub>n</sub> instability across successive generations is currently thought to involve the generation of abnormal DNA structures in replication slippage, as well as DNA repair and recombination, acting either separately or in combination.<sup>35</sup>

### Mechanisms of *ATXN3* mutation

The *ATXN3* gene is present in a large range of living organisms, but the expansion seems to have occurred only in the human lineage, after the divergence from other hominids. A recent study in human *ATXN3* has shown a bimodal distribution of (CAG)<sub>n</sub> alleles for all most frequent stable SNP-defined haplotypes. When flanking STR diversity was compared between modal (CAG)<sub>n</sub> alleles within each lineage, little differences have been noticed while alleles one or two repeats apart showed much higher genetic distances. Based on these results, a multistep mutation mechanism was suggested for the evolution of this *locus* in humans.<sup>36</sup> Little is known, however,



about the mutation process underlying the evolution and instability over mammalian and even primate evolution that led to the currently observed (CAG)<sub>n</sub> tract configuration and level of polymorphism. In *Pan troglodytes*, *Gorilla gorilla*, *Mus musculus* and *Gallus gallus*, the repetitive tract in the *ATXN3* gene is conserved, but shorter than in the human homologue.<sup>37</sup>

### **Ataxin-3 protein and its physiologic role**

The *ATXN3* gene encodes for ataxin-3, a protein with an approximate molecular weight of 42 kDa (in its longest form) that is ubiquitously expressed in neuronal and non-neuronal tissues.<sup>25, 38</sup> In the human brain, ataxin-3 is widely expressed in different regions, with variable expression levels. Granular layer and purkinje cells of cerebellum, hippocampus, striatum and pyramidal cells of motor cortex seems to be the regions of higher expression levels, whereas mesencephalon (substantia nigra pars compacta), occipital cortex, globus pallidus (internal and external) and white matter of cerebellum present lower expression patterns.<sup>25, 38</sup> In the *Drosophila* model, increased levels of the mutant *ATXN3* expression have been shown to underlie more severe degeneration and earlier onset of protein accumulation, suggesting that abnormal accumulation of the mutant protein is central to disease and degeneration.<sup>39</sup>

At the cellular level, *ATXN3* was found to be present in the cytoplasm (mitochondria included) and nucleus, with varying degrees of predominance depending on the cell type.<sup>40, 41</sup> In human brain cells, *ATXN3* is present mainly in the perikarya, but it was also detected on proximal processes, axons and nuclei. This heterogeneity suggests that the regulation of *ATXN3* expression levels and localization may have functional relevance.<sup>38</sup>

*ATXN3* participates in cellular protein quality control pathways but its biological function has not yet been completely understood. Studies in knockout mice for *ATXN3* have demonstrated a cytoplasmic increase in ubiquitinated proteins, supporting an *in vivo* role of this protein in the ubiquitin/proteasome pathway as a deubiquitinating enzyme (DUB).<sup>42, 43</sup> The covalent attachment of ubiquitin to a protein is a reversible signal that can alter a protein's function, control its trafficking, or mark it for degradation. Once attached, ubiquitin can be removed by DUBs, thereby generating a dynamic balance in ubiquitin signaling pathways.<sup>44</sup> Ataxin-3 can also regulate its own cellular turnover in a ubiquitin-dependent manner; abolition of its catalytic activity disrupts this regulation.<sup>45</sup> Hence, it is clear that ataxin-3 DUB activity is physiologically relevant. In addition, MJD mouse models have shown transcriptional deregulation, confirming a possible important role of *ATXN3* in transcription.<sup>46</sup>

ATXN3 has a papain-like fold and is essentially composed by a structured globular N-terminal domain, followed by a flexible unfolded C-terminal tail. The N-terminal domain, designed Josephin domain (JD), displays ubiquitin (Ub) protease activity, while the flexible tail presents two or three Ub-interacting motifs (UIMs), depending on the isoform, and a polyQ region of variable length.<sup>47, 48</sup> Notably, the most common isoform found in the human brain has three UIMs in the C-terminal region (Figure 4).<sup>49</sup> The highly conserved catalytic triad on the Josephin domain possesses the predicted catalytic aminoacids found in cysteine proteases: Cys14, His119 and Asn134.<sup>50</sup> Two other serine residues are important in the UIM regions for the interaction with Ub: Ser232 in UIM1 and Ser260 in UIM2.<sup>51</sup>

The Josephin domain is shared by three other human proteins (Ataxin-3 like, ATXN3L; Josephin-1, JOSD1; and Josephin-2, JOSD2), which together with ATXN3 form the Josephin family of DUBs.<sup>52</sup>

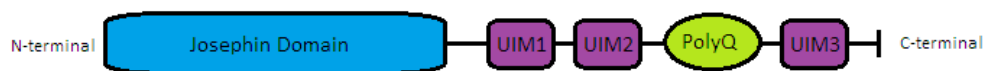


Figure 4 - Ataxin-3 protein structure. UIM – ubiquitin interacting motif; PolyQ – poly-glutamine tract.

A total of 56 human alternative splicing variants of *ATXN3* have been identified with, at least, 21 isoforms expected to be translated. The biological relevance of all these variants remains unknown.<sup>13, 53</sup>

### **Polyglutamine expansion of ATXN3 and neural cell death**

The expanded ATXN3 gains a neurotoxic function through yet unclear mechanisms. In brains of MJD patients, ATXN3 forms nuclear inclusions (NI) present only in neurons, however, more recently axonal inclusions have also been observed in fibers known to degenerate.<sup>41</sup> It is known, however, that this toxicity is linked to abnormal folding and aggregation of the mutant protein to itself as to other members of the cell quality control system, like proteasome constituents, ubiquitin and molecular chaperones.<sup>54, 55</sup> Additionally, intracellular inclusions have been linked to the disease pathogenesis as a major means of gain of toxicity by the expanded protein, through several possible mechanisms: a) hindrance of transcription, by direct decrease of gene expression regulation or through sequestration of other molecules involved in transcription regulation;<sup>56</sup> b) transcription alteration via formation of histone-deacetylating repressor complexes on target

chromatin regions;<sup>57, 58</sup> c) interference in the axonal transport, resulted from motor protein titration and physical blocking;<sup>59</sup> and d) other disturbances caused by the recruitment of Ub-binding proteins (since inclusions are heavily ubiquitinated) or other polyQ-containing proteins.<sup>60</sup> In contrast, NI have been recently suggested to play a protective role on MJD by sequestering the mutant proteins from the toxic interactions.<sup>61</sup>

Although the polyQ expansion appears to be the triggering factor leading to the development of MJD, ATXN3 regions outside the polyQ, as well as other protein properties seem also to define the development of MJD and its particular aspects. In animal and cellular models, the C-terminal fragment of expanded ataxin-3 alone was shown to be more toxic than the full-length protein, suggesting that ataxin-3 cleavage might be a contributor factor for the MJD pathology.<sup>62-64</sup> Indeed, the ataxin-3 expanded polyQ tract was shown to exhibit several potential cleavage sites for caspases.<sup>65</sup> Another study has also demonstrated a possible involvement of the N-terminal on MJD pathogenesis, since a mutant mouse model with an ATXN3 truncated form, containing only its N-terminal, presented typical neurological symptoms.<sup>66</sup>

Therefore, clinical variability of MJD is only partially explained by the size of the (CAG)<sub>n</sub> tract and by the ATXN3 protein itself, which leaves a residual variance that should be explained by other unknown factors. In addition to the wild-type protein, the three other JD-like containing proteins may exert similar functions to ATXN3 and compensate for its absence in knockout models.<sup>52, 67</sup> In fact, several studies on animal models have shown the importance of the normal non-mutated protein in the MJD pathogenesis. An interesting feature of non-expanded ataxin-3 is that it is also recruited to NI in several PoliQ diseases.<sup>68</sup> In the *Drosophila* model, for example, normal ATXN3 has a neurodegeneration repression effect *in vivo* (a protective role) by retarding and reducing the accumulation of the pathogenic protein. This suppressor activity requires ubiquitin-associated activities at the normal protein, since, by mutating the two most important serine residues in UIM1 and UIM2 (Ser232 and Ser260, respectively), the ability of ataxin-3 to suppress degeneration become compromised.<sup>69</sup>

### **1. 3 – Studies on the relevance of *ataxin* paralogues in SCAs**

Studies focusing on the modulation of parental genes by the respective expressed paralogues may be important to gain insight into the mechanisms by which *ATXN3* retrocopies may be functional relevant. The gene responsible for spinocerebellar ataxia type 1 (SCA1), *ATXN1*, has an

evolutionary conserved paralogue called *ATXN1-like*. Recent studies in fly and mouse models have shown that the overexpression of *ATXN1L* partially suppresses the neuropathology caused by the polyglutamine-expanded *ATXN1*, by inducing sequestration of polyglutamine-expanded *ATXN1* into nuclear inclusions.<sup>70, 71</sup> Similarly to MJD, molecular and genetic data suggested that SCA1 pathogenesis is caused by a gain of function mechanism, while other studies have shown, however, that the deletion of wild-type protein enhances disease pathogenesis. In SCA1, additionally, it has been suggested that increased levels of the *ATXN1L* retrogene ameliorate the clinical phenotype. These data indicated that both gain and partial loss of function may contribute to the disease progress, with the partial loss of *ATXN1* alone being sufficient to cause some transcriptional changes that are pathogenic in the cerebellum.<sup>70</sup>

The mechanism underlying the functional relevance of *ATXN1L* may be related to the shared protein interactions by *ATXN1* and *ATXN1L*. The two proteins were shown to share the transcriptional regulator Capicua (CIC) that is known to form complexes (*Atxn1-CIC* and *Atxn1L-CIC*) that bind the promoters of target genes, repressing them effectively. Thus, in *ATXN1<sup>-/-</sup>* mice, the overexpression of *ATXN1L* rescues the level of stable complexes with CIC, maintaining its function as a transcription repressor. This mechanism is an evidence that transcribed paralogues can, in some degree, replace the loss-of-function observed in SCA1.<sup>70</sup> Additionally, it was proposed that, as *ATXN1L* compete with wild-type and mutant *ATXN1* for association with CIC, the increased levels of free mutant *ATXN1* lead to an increase of aggregation and nuclear inclusions.<sup>70, 71</sup> This data provide genetic evidence that evolutionary conserved paralogues may have an important role in the mechanisms of pathogenesis of neurodegenerative diseases.

As for *ATXN3* and MJD/SCA3, Weeks (2011) performed an in vitro study of the human *ATXN3L* retrocopy (here named *ATXN3L1*) by analysing the crystal structure of its predicted protein.<sup>52</sup> The authors found that although ataxin-3 and *ATXN3L* adopt similar folds, they bind ubiquitin in different, overlapping sites. Additionally, by mutating ataxin-3 at selected positions (and introducing the corresponding *ATXN3L* residue), only three mutations were sufficient to increase the catalytic activity of ataxin-3. This suggested that *ATXN3L* Josephin domain could be significantly more efficient than the ataxin-3 domain itself, opening a broad road for further research concerning the study of this *ATXN3* paralogue.

# Aims

---

Two human copies of the *ATXN3* gene, formed by retrotransposition events, were recently discovered and annotated in databases. However, they remain unexplored as their description has not been received much attention. Therefore, we wanted to better characterize these retrocopies since, if transcribed, they may play a role in the pathogenesis of MJD. Additionally, by comparing the (CAG)<sub>n</sub> tracts among the 3 paralogue sequences, we hope to gain insight into the repeat expansion mechanism that occurred at *ATXN3* in the human lineage.

Therefore, the aims of this project were:

- 1) to estimate the onset of the retrotransposition events that originated the two human *ATXN3* paralogues;
- 2) to compare the rates of evolution and selective constrains underlying the three *ATXN3* paralogues throughout the primate lineage;
- 3) to gain insight into the mechanisms that led to the human-specific (CAG)<sub>n</sub> expansion and tract configuration of *ATXN3*; and
- 4) to investigate if the human *ATXN3* paralogues are transcribed and, if so, their mRNA transcription pattern.



## Chapter 2

# Material and methods

---





## Chapter 2

# Material and methods

---

To achieve the aims of this project, an evolutionary approach (Part 1) was firstly performed, followed by functional assays (Part 2). We have done evolutionary studies to compare the rates of evolution and selective constraints among *ATXN3*, *ATXN3L1* and *ATXN3L2* genes, and to date the onset of the retrotransposition events. We started by aligning homologous sequences from several primate species (followed by more distant mammals) available in databases. Afterwards, we calculated genetic distances and determined both the origin of *ATXN3* retrocopies and the underlying selective pressures throughout evolution. In Part 2, we started by comparing the Open Reading Frames (ORFs) among *ataxin-3* paralogues. For *ATXN3L1*, which showed a conserved ORF, we assessed the mRNA expression pattern in different human tissues, by Reverse Transcriptase-PCR (RT-PCR). The two approaches, Part 1 and Part 2, are described below in detail. As it was previously referred, *ATXN3* paralogues, described as *ATXN3L* and *LOC100132280*, will be named as *ataxin-3 like 1 (ATXN3L1)* and *ataxin-3 like 2 (ATXN3L2)*, respectively, for easier reading.

### Subjects

We analysed a total of 49 healthy human subjects from European (n=6) and Asian (n=43) origins. DNA samples were available at our lab, previously coded to ensure confidentiality. Informed consent was provided by all individuals.

As for non-human primate species, we analysed DNA samples from *Pan troglodytes* (chimpanzee) (n=2), *Gorilla gorilla* (gorilla) (n=3), *Pongo abelii* (orangutan) (n=2), *Macaca mulatta* (rhesus monkey) (n=5), *Macaca fascicularis* (cynomolgus monkey) (n=3), *Callithrix jacchus* (marmoset) (n=3) and *Papio* (baboon) (n=3). DNA was quantified with NanoDrop spectrophotometer (Thermo Scientific) to further make work aliquots with a DNA concentration of approximately 30 ng/ $\mu$ L.

## Part 1 – Evolutionary approach

### Compilation and alignment of *ATXN3L1* and *ATXN3L2* sequences

To study the origin of the retrotransposition events, we first needed to identify *ATXN3* paralogues in public databases. We started by searching for annotated *ATXN3L1* and *ATXN3L2* orthologues in primates and other mammalian species. Sequences of *ATXN3L1* were available in Ensembl database (<http://www.ensembl.org/>) for *Homo sapiens*, *Pan troglodytes*, *Pongo abelii*, *Macaca mulatta*, *Callithrix jacchus* and *Nomascus leucogenys*. The sequence from human *ATXN3L2*, termed LOC100132280, was obtained from the National Center for Biotechnology Information (NCBI) database (<http://www.ncbi.nlm.nih.gov/>), but no other annotated orthologues were found in either Ensembl or NCBI Gene databases. Therefore, we performed BLAT and BLAST algorithms provided by University of California, Santa Cruz (UCSC) Genome Browser (<http://genome.ucsc.edu/>) and NCBI databases, respectively. This way, we obtained not yet annotated homologous sequences of *ATXN3L2* in the genome of several primate species. Following the same strategy, we searched for additional *ATXN3L1* sequences of primate species not retrieved from Ensembl. The same approach was also applied to search for *ATXN3L1* and *ATXN3L2* orthologues in other mammals; in this case, however, we used mostly the Trace Archives specialized search of BLAST algorithm (NCBI).

The alignment of the collected sequences (coding regions for *ATXN3*) was performed in Geneious Pro, 5.5.6 software and the homology (percentages of pairwise identity and identical sites) between sequences was obtained among paralogues of primate species. This allowed to discern whether two independent retrotransposition events occurred or, alternatively, one of the copies duplicated from the most ancient retrocopy.

Next, to identify the transcript(s) of the parental gene involved in the retrotransposition events, we aligned *ATXN3L1* and *ATXN3L2* with the 21 protein coding transcripts described for *ATXN3*. Transcripts were retrieved from Ensembl: *ATXN3-001*, *ATXN3-003*, *ATXN3-004*, *ATXN3-005*, *ATXN3-008*, *ATXN3-015*, *ATXN3-017*, *ATXN3-019*, *ATXN3-020*, *ATXN3-026*, *ATXN3-029*, *ATXN3-032*, *ATXN3-201*, *ATXN3-202*, *ATXN3-203*, *ATXN3-204*, *ATXN3-205*, *ATXN3-206*, *ATXN3-207*, *ATXN3-208* and *ATXN3-209*. In humans, *ATXN3-001* has been described as the most common transcript, followed by *ATXN3-003*, *004* and *005*.

### **Phylogenetic trees and genetic distances**

Based on the aligned sequences, we calculated genetic distances among orthologues and paralogues, and constructed phylogenetic trees for all species, by using the same software, Geneious Pro 5.5.6. To build the trees, we applied the Neighbor-Joining method and the Tamura-Nei genetic distance model (chicken, *Gallus gallus*, was included as outgroup). This approach was performed to identify *ATXN3L1* and *ATXN3L2* orthologues. Since the Josephin domain is the most conserved region of *ATXN3* and the (CAG)<sub>n</sub> region is usually a source of much variation, we assessed genetic distances and built phylogenetic trees for all three paralogues by using: 1) the entire sequences; 2) the sequences without the (CAG)<sub>n</sub> tract; and 3) the Josephin domain alone.

### **Synten of *ATXN3*, *ATXN3L1* and *ATXN3L2***

We analysed the genes flanking *ATXN3*, *ATXN3L1* and *ATXN3L2* in a region of, at least, 600000 base pairs (bp) to elucidate if these regions were conserved along the primate lineage, confirming thus gene orthologies. This was performed by using NCBI Gene database, UCGC Genome Browser genomes and Ensembl Synteny analyses. If two flanking regions of a given gene have obvious collinear genomic features, they are syntenous. From this information, therefore, we can conclude that predicted orthologue genes having syntenous flanking regions have higher probability to be real orthologues. In addition, the same approach was applied to analyse sequences of non-primate mammals: their flanking regions were compared to those retrieved from primates in order to identify if they were ancestral retrocopies of *ATXN3L1* and/or *ATXN3L2*, or alternatively, originated from independent retrotransposition events.

### **dN/dS ratio (omega) calculations**

Evolutionary pressures on proteins can be quantified by the ratio of substitution rates at non-synonymous and synonymous sites. We calculated dN/dS values (omega) of *ATXN3* and *ATXN3L1* genes for different primate species, using the DnaSP v5 software. The dN value represents the rate of non-synonymous substitutions (nucleotide changes that result in amino acid alterations) per site, whereas dS is the rate of synonymous substitutions (silent alterations that do not change the amino acid residue). Thus, a dN/dS ratio (omega) is expected to exceed unity if substitutions are equally frequent at all three codon positions, indicating that no constraints have been underlying the gene. On the other hand, omega less than unity is expected if selection suppresses protein changes, suggesting that the gene function has been kept in check. These calculations

allowed us to compare the selective constrains that have been underlying *ATXN3* paralogues in several lineages of primates. As *ATXN3L2* presented a disrupted ORF, with the start codon mutated and several premature stop codons, we did not proceed with this calculation for this retrocopy.

### Primer design, DNA amplification and sequencing

Based on the alignment of gene orthologues, primers were designed to specifically amplify *ATXN3* (exonic regions), *ATXN3L1* and *ATXN3L2* in all primate species. Exons of *ATXN3* annotated in NCBI database represent only those included in the most frequent transcript (11 exons). When analysing all *ATXN3* protein coding transcripts represented in Ensembl data base, we retrieved other exons transcribed only in some transcripts (protein coding) and others only predicted to be transcribed. As reference, in this study, we used all *ATXN3* exons annotated in Ensemble (21 in total), although for primer design, we selected only exons known to be transcribed. Figure 5 summarizes this information.

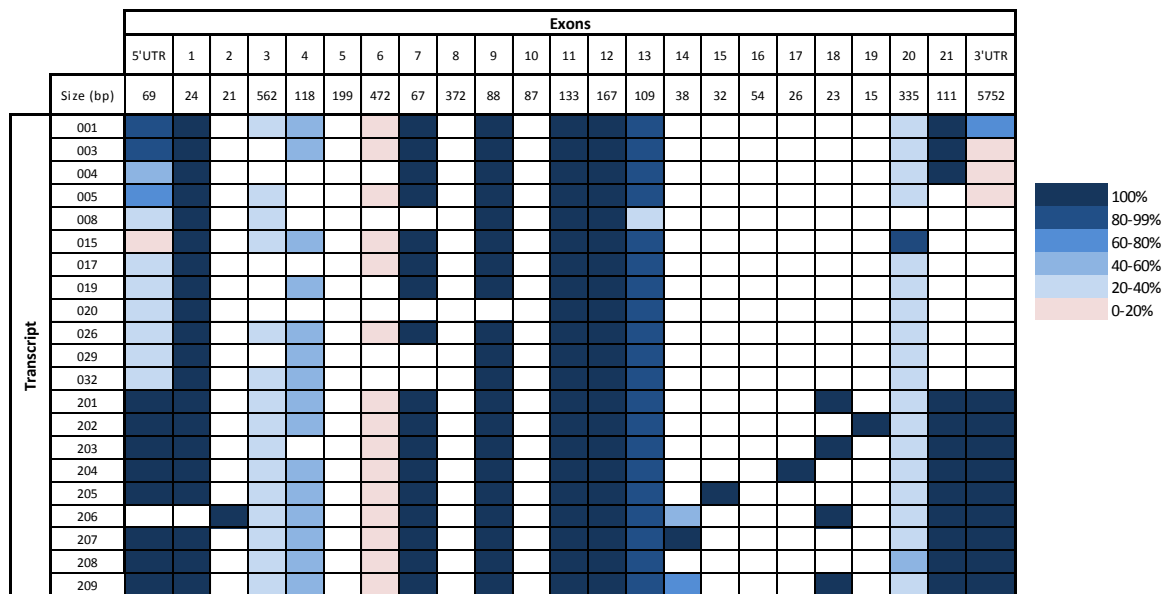


Figure 5 – Exons of the *ATXN3* protein coding transcripts, according to Ensemble database. As some transcripts do not contain the entire exon sequence, percentages of coding sequences relatively to the total number of nucleotides of the exon were calculated and represented with different colors. Color legend is on the right side of the figure.

For primate species that lacked the annotated sequence of one or more *ATXN3* paralogues (namely chimp, gorilla, orangutan, gibbon, marmoset, cynomolgus monkey and baboon), we performed sequencing; the same PCR conditions and oligo pairs were used, regarding the sequence conservation of the selected primers for most of the primates.

For primer design, we selected the following criteria: 20 to 24 base pairs; a percentage of GC between 45 and 60%; and a melting temperature comprised between 56 and 64°C. The OligoCalc algorithm (<http://www.basic.northwestern.edu/biotools/oligocalc.html>) was used to calculate melting temperatures and GC percentages, as well as to predict the formation of hairpins and primer dimers. Primers predicted to form these structures were rejected. In addition, to avoid the formation of primer-dimers, we used the AutoDimmer software to select primers with no such predicted structures. Finally, to analyze primers' specificity for each species, we performed an UCSC Genome Browser *In-Silico PCR* (<http://genome.ucsc.edu/cgi-bin/hgPcr?command=start>). The list of selected primers is shown in Tables 1 and 2.

**Table 1 – Primers designed to specifically amplify *ATXN3* exons.**

<b><i>ATXN3</i></b>					
Reaction	T <sub>annealing</sub>	Exons	Amplicon	Primer sequence	Species
Pentaplex	57°C	Part of 2 and 3	F2-R3	F2 – GTGGTAAGCTGAGATTGCTCC R3 – CCAGCTGATGTGCAATTGAGG	<i>H. sapiens</i> <i>P. troglodytes</i> <i>G. gorilla</i> <i>P. abelii</i>
		Part of 6	F6-R6	F6 – CACAACAAACATAGCTACACTTCC R6 – AAGGCTACAGGGCAGATGCT	
		9	F9-R9	F9 – CCTGGCCAATGTGGCAAATG R9 – CACTGTCATCTAATGTGCTG	
		12	F12-R12	F12 – GGTTGCAGTTATTACCAAGTGC R12 – GAAATCTAAAGGAAAGCCAC	
		16	F16-R16	F16 – CTGATCCTAGGTGAGAAACAGA R16 – GGCCGTGTGCTAGTATTGTTG	
Hexaplex	60°C	4	F4-R4	F4 – GCACGCTAATGACAGTTTGTATCC R4 – GGTGAAACCCCACTATCTCTAC	<i>H. sapiens</i> <i>P. troglodytes</i> <i>G. gorilla</i> <i>P. abelii</i>
		Part of 6 and 7	F6'-R7	F6' – CAACAGTCCAGAGTATCAGAGGC R7 – GACAGGACCTCCCTTTGTTGCC	
		11	F11-R11	F11 – TCCAGTGTCTGTGCTGCCTTTT R11 – AGTCGCCAACAACAAGGACC	
		17	F17-R17	F17 – GGAAAGGCATCTCTGGGGAG R17 – GAAGTTTGACACGAGCCTGGAC	
		18,19 and part of 20	F18-R20	F18 – CCACTCCTGGCCATGATAGGT R20 – GAATGGTGAGCAGGCCTTACC	
21	F21-R21	F21 – CTGGTGGCTATCTGGGATTAGGA R21 – GGACCCTATGCTGTAATCACACAG			
Duplex 1	64°C	1	F1-R1	F1 – CGTGTCCCCGGCGTTCACCTC R1 – AGATCGGCATGGGGCGACT	<i>H. sapiens</i> <i>P. troglodytes</i> <i>G. gorilla</i> <i>P. abelii</i>
		3	F3-R3'	F3 – GCAAGAAGGCTCACTTTGTGCTC R3' – CAAGGGTGGGGTGGGAAA	
Duplex 2	64°C	13 and 14	F13-R14	F13 – ACGCCCAGCCAGAAGAGTAG R14 – CTCCTGACCTCAGGCAATCTG	<i>H. sapiens</i> <i>P. troglodytes</i> <i>G. gorilla</i> <i>P. abelii</i>
		Part of 20	F20-R20	F20 – GGCCAGCCACCAGTTCAGGAG R20' – TCCTCTCTGCCTTGGTTTCCC	

Table 2 - Primers designed to specifically amplify *ATXN3L1* and *ATXN3L2* genes.

ATXN3L1 and ATXN3L2					
Reaction	T <sub>annealing</sub>	Gene	Amplicon	Primer sequence	Species
Singleplex 1	59°C	<i>ATXN3L1</i>	L1F1-L1R1	L1F1 – CTCTAACTAGGATACCAGCAAAG L1R1 – GGAAAAAGTTCTATGGCAAGAGC	<i>H. sapiens</i> <i>P. troglodytes</i> <i>G. gorilla</i> <i>P. abelii</i> <i>M. mulatta</i> <i>N. leucogenys</i> <i>C. jacchus</i>
Singleplex 2	57°C	<i>ATXN3L2</i>	L2F1-L2R1	L2F1 – CATTAACCAAAGAAAGTGGGATAC L2R1 – GGAATCCTATGCTGTAATCACAC	<i>H. sapiens</i> <i>P. troglodytes</i> <i>G. gorilla</i> <i>P. abelii</i> <i>M. mulatta</i> <i>N. leucogenys</i>
Singleplex 3	57°C	<i>ATXN3L2</i>	cL2F1-cL2R2	L2F1 – CATTAACCAAAGAAAGTGGGATAC L2R1 – GGAATCCTATGCTGTAATCACAC	<i>C. jacchus</i> *

\* As the *ATXN3L2* sequence of *Callithrix jacchus* differed more than 25% from the other sequences, different primers were designed for this species.

We optimized polymerase chain reactions (PCRs) with the designed primers described above. PCRs were done in a final volume of 10 µL, with 1x of Taq polymerase 2x (MyTaq; BIOLINE) and 0.25 µM (in singleplex reactions) or 0.125 µM (in multiplex reactions) of each primer. PCR conditions are described in Figure 6.

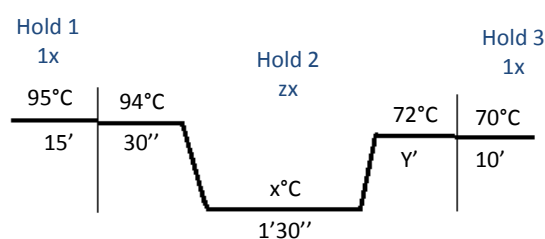


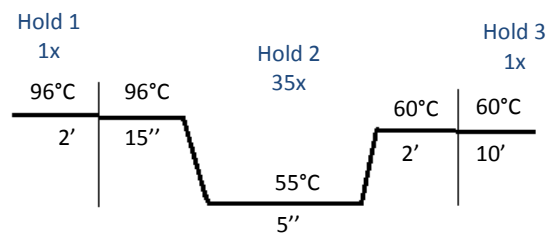
Figure 6 – General PCR protocol with time (‘ - minutes; ‘‘ - seconds) and temperatures (°C) for the amplification of *ATXN3*, *ATXN3L1* and *ATXN3L2* loci. X, Y and Z vary according to the multiplex/singleplex reaction: X (annealing temperature) – described in Tables 1 and 2; Y (extension time) – 2 minutes for singleplex, and 1 minute for multiplexes; Z (number of cycles) – 30 for *ATXN3*; 35 for *ATXN3L1*; and 40 for *ATXN3L2*.

To confirm specificity of DNA amplification, PCR products were submitted to electrophoresis in a polyacrylamide gel prepared with 0.375M Tris/HCl gel buffer (pH=8.8). The gel was obtained by mixing 3 mL T<sub>9</sub>C<sub>5</sub> (9% acrylamide, 5% N-N-metileno-bis-acrylamide), 170 µL of 1% ammonium persulfate and 7 µL of TEMED as catalysis agent. Glass supports, with one side covered by a hydrophilic gel-bond film, were used to obtain a 3 mm thick gel. After loading the samples in the

gel, two paper strips soaked in buffer were used at both anode and cathode to allow the horizontal run. To monitor the process of electrophoresis, we added bromophenol blue dye to the anode strip. The electrophoretic system was submitted to refrigeration at 4°C, and to a voltage between 220 and 250 V. The gel was then submitted to silver staining. The coloration method comprised (1) a fixation step of the DNA, with 10% ethanol for 10 minutes followed by 1% nitric acid for 5 minutes; (2) two washes with deionized water, of about 20 seconds each; (3) a coloration step with 0,2% silver nitrate solution, for 20 minutes; (4) two washes again with deionized water, for 20 seconds each; and finally, (5) a revelation step of the DNA fragments with a solution of 0,28 M sodium carbonate and 0,02% formaldehyde. The revelation reaction was stopped with 10% acetic acid for approximately 30 seconds. The resulting gels were washed with water and dried at room temperature.

After optimizing PCR reactions, we proceeded with sequencing. Products were first purified with ExoStar (GE Health Care), for 15 minutes, at 37°C, followed by 15 minutes at 80°C to inactivate the enzyme. ExoStar contains a mix of Alkaline Phosphatase and Exonuclease 1, formulated to remove unincorporated primers and nucleotides from the resulting PCR products.

For sequencing procedures, we used Big Dye Terminator Cycle Sequencing Ready Reaction kit (Applied Biosystems) as mix containing normal deoxynucleotides, dye dideoxynucleotides, buffer and AmpliTaq DNA Polymerase. For each reaction of 5 µL, we combined 1 µl of Big Dye mix, sequencing buffer 2.5x, and 0.5 µM of primer with 2.5 µL of purified PCR product. Cycling conditions are described in Figure 7.



**Figure 7 - General protocol with time ('- minutes; '' – seconds) and temperatures (°C) for sequencing ATXN3, ATXN3L1 and ATXN3L2 PCR products**

Finally, sequencing products were purified using Sephadex columns, a cross-linked dextran-gel used to separate the low (unincorporated nucleotides and primers) from high molecular weight molecules of DNA (sequencing products); formamide was then added to the final product to increase the stability of single-stranded DNA for the capillary electrophoresis run in an ABI 3130

Genetic Analyzer. Sequences were analysed with Sequencing Analysis v5.2 software (Applied Biosystems).

### **Nucleotide diversity of human *ATXN3* retrocopies**

To evaluate the nucleotide diversity underlying *ATXN3L1* and *ATXN3L2* human retrocopies, we aligned sequences obtained from the previous step and annotated all variations relative to the reference sequence. We calculated frequencies of variants and predicted the alterations in the respective putative coding sequences.

### **(CAG)<sub>n</sub> tract analysis in *ATXN3*, *ATXN3L1* and *ATXN3L2***

We compared the (CAG)<sub>n</sub> tract configuration among the three paralogues for all species based on (1) sequences collected from NCBI, Ensembl and UCSC Genome browser, and (2) results obtained from sequencing (*Pan troglodytes*, *Gorilla gorilla*, *Pongo abelii*, *Macaca mulatta*, *Macaca fascicularis* and *Papio*).

## **Part 2 – Functional approach**

### **ORF prediction for *ATXN3L1* and *ATXN3L2***

Since little was known regarding the functional relevance of *ATXN3L1* and *ATXN3L2*, we started by predicting the respective ORFs. Predictions were performed by using NCBI Open Reading Frame Finder (ORF Finder - <http://www.ncbi.nlm.nih.gov/projects/gorf/>) graphical analysis tool.

### **Conservation of the *ATXN3L1* putative protein**

The expected protein coding sequence of *ATXN3L1* was aligned and compared to the most common coding sequence of the parental gene, the *ATXN3-001* transcript. The Josephin domain, UIMs and polyQ tract were annotated in sequences, as well as other important sites/amino acids such as the nuclear export signals NES77 and NES141, the nuclear localization signal NLS273, and the catalytic amino acid triad. With this comparative study, we were able to determine if *ATXN3L1* conserved the most important functional components of the parental protein and, thus, a



potential similar functional activity. As ATXNL2 presented premature stop codons and is unlikely to be translated, we did not proceed with further analysis in this retrocopy.

### **mRNA expression of *ATXN3L1***

To evaluate if *ATXN3L1* is transcribed, primers were designed to specifically amplify ATXN3L1 cDNA. The cDNA of 16 different human tissues was previously obtained by conversion from mRNA, using reverse transcriptase PCR (RT-PCR). Analysed tissues included: ovary, bladder, trachea, esophagus, thymus, thyroid, colon, kidney, skeletal muscle, testis, small intestine, heart, spleen, prostate, liver and brain.

We started by aligning the *ATXN3L1* sequence with the concatenated exons of parental *ATXN3*; this way, we annotated the exons' limits of the parental gene in the retrocopy. Then, primers were designed in regions flanking the exon junctions of the corresponding *ATXN3* exons. Primers are listed in Table 3.

**Table 3 - Primers designed to specifically amplify ATXN3L1 in human cDNA.**

<b>cDNA <i>ATXN3L1</i> primers</b>				
<b>Reactions</b>	<b>T<sub>annealing</sub> (°C)</b>	<b><i>ATXN3</i> exons</b>	<b>Amplicon</b>	<b>Primer sequence</b>
ATXN3L1 cDNA	58	UTR, 1 and 2	UTR-E2F - UTR-E2R	UTR-E2F - GCATACAACATCTCCGGCATACC UTR-E2R - CAGACAGTGCTGAGCACACAGGA
		7 and 8	E7-8F - E7-8R	E7-8F - CAGTGTCGAAGAGATGGATAC E7-8R - CCTCATCTGGTCTGATGTTCCAGACT
		8 and 9	E8-9F - E8-9R	E8-9F - GAACTAAGCCGCAAGAAACC E8-9R - GCAGGAGTTACACATGATGTCTTTGGA
		10 and 11	E10-11F - E10-11R	E10-11F - GGGCCACAGTTCATACCTACAC E10-11R - TGTCGACAGCGCCTGACTG

PCR conditions were similar to those used for sequencing described above (Figure 1), using 58°C as annealing temperature, 1 minute for extension time, and 35 cycles. To confirm amplification, we have done an electrophoresis run, followed by silver stained, as described above. The expected product sizes are approximately 110 bp for UTR-E2, 140 bp for E7-8, 130 bp for E8-9 and 110 bp for E10-11 amplicon.



## Chapter 3

# Results

---



## Chapter 3

# Results

---

## Part 1 – Evolutionary history of *ATXN3* paralogues

### 1. Onset of *ATXN3* retrocopies

#### 1.1. Identification of *ATXN3L1* and *ATXN3L2* orthologues in the primate lineage

We retrieved the coding sequence of *ATXN3* for several primates (Table 4). Although the size and location of the collected sequences varied among species, we found homogeneity on both parameters for great apes: humans, chimpanzees, gorillas and orangutans. Notably, when compared to the human sequence, differences were found in *ATXN3* sequences of 1) gorilla, in which exon 2 is not annotated as being translated; 2) orangutang, which does not possess a complete annotated transcription for exon 2; and 3) tarsier, for which databases do not display the complete sequence.

As for *ATXN3* retrocopies, a high homology to the human *ATXN3L1* and *ATXN3L2* sequences was found in *Pan troglodytes*, *Gorilla gorilla*, *Pongo abelii*, *Macaca mulatta*, *Callithrix jacchus*, *Nomascus leucogenys*, *Tarsius syrichta* and *Otolemur garnettii*. Based on the percentages of homology to human retrocopy sequences, calculated through NCBI-BLAST analysis, we designated the collected sequences as *ATXN3L1* or *ATXN3L2* (Tables 5 and 6). Tarsier and bushbaby sequences (found to be aligned with the human *ATXN3L1*, by using the “comparative genomics” tool of Ensembl) presented a high divergence from both human *ATXN3L1* and *ATXN3L2*. For this reason, the alignment performed by Ensembl did not give us a reliable hint on the origin of these copies in the tarsier and bushbaby species; we, thus, named them as “*ATXN3L?*” at this point of the project. One other sequence for tarsier was retrieved from trace archives specialized BLAST of NCBI (Table 7). No sequence homologies were found in other primate species with genomes included in trace archives, such as *Pan paniscus* (bonobo), *Pongo pygmaeus* (bornean orangutan), *Lemur Catta* (ring-tailed lemur), *Aotus nancymae* (Nancy Ma's night monkey), *Callicebus moloch* (red-bellied Titi), *Chlorocebus aethiops* (grivet or green monkey), *Hylobates concolor* (black

crested gibbon), *Ateles geoffroyi* (Geoffroy's spider monkey), and others. This might be partially explained by the fact that trace archives contain highly fragmented sequences and that these data files are not official entries of the GenBank database, thus, with no associated feature annotations.

Additionally to the *ATXN3L1* and *ATXN3L2* sequences retrieved above, other less homologous sequences (displaying, even though, more than 80% sequence identity to human *ATXN3* coding sequence) were found in Ensembl data base for *Macaca mulatta* and *Callithrix jacchus* (Table 7). These sequences, here classified also as *ATXN3L?*, were aligned with all *ATXN3*, *ATXN3L1* and *ATXN3L2* (Figure A1 of the Appendix section) to next calculate genetic distances among them in topic 1.3 of the Results section.

Table 4 – *ATXN3* genomic and coding sequences compiled for 9 primate species. a) Genomic location of *Macaca mulatta ATXN3* was assessed through the UCSC Genome Browser; b) Incomplete annotated sequence.

### ATXN3

Species	Source	Location	Size of genomic/coding sequences (bp)
<i>Homo sapiens</i>	Ensembl	chr14: 92,524,896-92,572,965	48070/1113
<i>Pan troglodytes</i>	Ensembl	chr14: 91,609,895-91,656,494	46600/1086
<i>Gorilla gorilla</i>	Ensembl	chr14: 73834975-73877125:-1	42151/937
<i>Pongo abelii</i>	Ensembl	chr14: 93360787-93404545:-1	43759/996
<i>Nomascus leucogenys</i>	Ensembl	GL397280.1: 29388208-29433598:-1	45391/1080
<i>Macaca mulatta</i>	Ensembl	Scaffold 1099553000000: 10,051-45,764 (Chr7) <sup>a</sup>	35714/1059
<i>Callithrix jacchus</i>	Ensembl	chr10: 117,483,688-117,689,710	206023/1101
<i>Otolemur garnettii</i>	Ensembl	GL873539.1: 18068653-18100321:-1	31669/867 <sup>b</sup>
<i>Tarsius syrichta</i>	Ensembl	GeneScaffold_567: 5192-20314:-1	15123/1089 <sup>b</sup>

Table 5 – *ATXN3L1* sequences compiled for 7 primate species.

### ATXN3L1

Species	Source	Location	Size (bp)
<i>Homo sapiens</i>	Ensembl	chrX: 13336770-13338518:-1	1068
<i>Pan troglodytes</i>	Ensembl	chrX: 13251386-13252912:-1	1062
<i>Gorilla gorilla</i>	UCSC	chrX: 13240248-13241315	1068
<i>Pongo abelii</i>	Ensembl	chrX: 13230000-13231067:-1	1069
<i>Nomascus leucogenys</i>	Ensembl	GL397281.1: 10,640,111-10,641,175	1065
<i>Macaca mulatta</i>	Ensembl	chrX: 11,016,157-11,017,194	1065
<i>Callithrix jacchus</i>	Ensembl	chrX: 11,295,476-13,523,142	1062

**Table 6 - ATXN3L2 sequences compiled for 7 primate species, mainly by using BLAST and BLAT algorithms, from NCBI and UCSC Genome Browser, respectively.**

**ATXN3L2**

Species	Source	Location	Size (bp)
<i>Homo sapiens</i>	NCBI (Gene ID: 100132280)	8q23.2	1109
<i>Pan troglodytes</i>	BLAST/BLAT search	chr8: 109011629-109012743	1115
<i>Gorilla gorilla</i>	BLAST/BLAT search	chr8: 109925669-109926225	1084
<i>Pongo abelii</i>	BLAST/BLAT search	chr8: 117579303-117579866	1075
<i>Nomascus leucogenys</i>	BLAST/BLAT search	GL397267:24598565-24599123	1075
<i>Macaca mulatta</i>	BLAST/BLAT search	chr8:112963801-112964361	1079
<i>Callithrix jacchus</i>	BLAST/BLAT search	chr19:6218140-6218664	1056

Table 7 - Additional ATXN3 paralogues found in primates.

**ATXN3L?**

Species	Source	Location/Accession
<i>Callithrix jacchus</i>	Ensembl	chr20: 28,630,500-28,631,553
<i>Macaca mulatta</i>	Ensembl	chr11: 122178105-122178311
<i>Tarsius syrichta</i> (seq1)	Ensembl - Comp. genomics	Unknown
<i>Tarsius syrichta</i> (seq2)	BLAST - Trace Archives	Unknown
<i>Otolemur garnettii</i>	Ensembl - Comp. genomics	Unknown

Next, to undoubtedly classify *ATXN3L1* and *ATXN3L2* orthologues in primates and distinguish them from other independent-origin copies, we started by assessing the synteny of all *ATXN3*, *ATXN3L1* and *ATXN3L2* (Table 8 – A, B and C). For *Tarsius syrichta* and *Otolemur garnettii*, this search over flanking regions was not possible in any database. To assess their origin, we compared genetic distances and constructed phylogenetic trees in the next topic (1.2 of the Results section).

By analysing the synteny of *ATXN3* genomic regions, we noticed that up and downstream genes were constant throughout the primate lineage. Upstream genes include, mainly, *CPSF2* (cleavage and polyadenylation specific factor 2) and *NDUFB1* (NADH dehydrogenase 1 beta subcomplex), whereas downstream genes include *TRIP11* (thyroid hormone receptor interactor 11) and *FBLN 5* (fibulin 5). The sodium/potassium/calcium exchanger gene was only detected upstream to *ATXN3* of some primates (*Homo sapiens*, *Pongo abelii* and *Otolemur garnettii*); however, it may be also present in other species since some genes may have currently not been annotated for all species. The same applies to *PTMAP7* (prothymosin, alpha pseudogene 7), only detected downstream to *ATXN3* in humans and chimps.

When we analysed the synteny of *ATXN3L1*, we have found a conserved upstream region across the 7 primates, which included *TCEANC* (*transcription elongation factor A N-terminal and central domain containing*) and *EGFL6* (*epidermal growth factor-like protein 6*) genes. In addition, downstream sequences contained typically two or three of the following genes: *FAM9C* (*family with sequence similarity 9, member C*), *TMSB4X* (*thymosin beta 4, X-linked*) and *TLR8* (*toll receptor 9 precursor, CD289 antigen*).

As for *ATXN3L2*, the majority of primate species presented *KCNV1* (*potassium channel, subfamily V, member 1*) in the 5' flanking region and *CSMD3* (*CUB and Sushi multiple domains 3*) downstream. Human 3' genes were distinct from all other primates: two pseudogenes were present, one of them, with the respective parental gene located upstream of *ATXN3* in humans (*NADH dehydrogenase 1 beta subcomplex pseudogene*). Upstream regions are, however, conserved, supporting this retrocopy is the human *ATXN2* orthologue. Therefore, from this analysis, we confirmed that the majority of collected sequences had been correctly classified based on sequence homology. In the case of the marmoset retrocopy, however, previously classified as *ATXN3L2*, synteny did not favour our first hypothesis of a common origin shared with other *ATXN3L2* orthologues. We, therefore, reclassified this *Callithrix jacchus* sequence as *ATXN3L?2* and performed additional analyses to identify its origin (*ATXN3L?1* denotes the first detected independent-origin retrocopy for *Callithrix jacchus*).



Table 8 - Synteny of *ATXN3* (A), *ATXN3L1* (B) and *ATXN3L2* (C) through evidence based on collinearity of genes in primates. Two or three genes were retrieved within an interval of 600000 bp, except those marked with an asterisk (\*) which were found outside this interval. a, b and c next to the species name represent the source of the information: Ensembl, NCBI and UCSC Genome Browser, respectively. When information for a specific gene or location was retrieved from a source different from the indicated at species name, a, b or c specify the source.

A - *ATXN3*

Species	Flanking 5'	Gene location	Flanking 3'
<i>Homo sapiens</i> <sup>a</sup>	NM_153648.3 (sodium/potassium/calcium exchanger 4 isoform 3) CPSF2 (cleavage and polyadenylation specific factor 2) NDUFB1 (NADH dehydrogenase (ubiquinone) 1 beta subcomplex)	Chr.14q32.1	PTMAP7 (prothymosin, alpha pseudogene 7) <sup>b</sup> TRPIP11 (thyroid hormone receptor interactor 11) FBLN5 (fibulin 5)
<i>Pan troglodytes</i> <sup>a</sup>	CPSF2 (cleavage and polyadenylation specific factor 2) NDUFB1 (NADH dehydrogenase (ubiquinone) 1 beta subcomplex)	Chr.14: 91,609,895-91,656,494	LOC100614355 (prothymosin alpha-like) <sup>b</sup> TRPIP11 (thyroid hormone receptor interactor 11) FBLN 5 (fibulin 5)
<i>Gorilla gorilla</i> <sup>a</sup>	CPSF2 (cleavage and polyadenylation specific factor 2) NDUFB1 (NADH dehydrogenase (ubiquinone) 1 beta subcomplex)	Chr.14: 73,834,975-73,877,125	TRIP11 (thyroid hormone receptor interactor 11) FBLN 5 (fibulin 5)
<i>Pongo abelii</i> <sup>a</sup>	LOC100443929 (cleavage and polyadenylation specific factor 2) <sup>b</sup> NDUFB1 (NADH dehydrogenase (ubiquinone) 1 beta subcomplex) LOC100446268 (Retinal Na+/Ca2+/K+ Exchanger)	Chr.14: 93,360,787-93,404,545	LOC100441607 (thyroid hormone receptor interactor 11) FBLN 5 (fibulin 5)
<i>Macaca mulatta</i>	?	Chr.7 <sup>b</sup> (Scaffold 1099553000000: 10,051-45,764) <sup>a</sup>	LOC100423180 (uncharacterized) <sup>b</sup>
<i>Callithrix jacchus</i> <sup>b</sup>	CPSF2 (cleavage and polyadenylation specific factor 2) LOC100414432 (uncharacterized)	Chr.10: 117,483,688-117,689,710 <sup>a</sup>	TRPIP11 (thyroid hormone receptor interactor 11) FBLN 5 (fibulin 5) <sup>a</sup>
<i>Nomascus leucogenys</i> <sup>a</sup>	CPSF2 (cleavage and polyadenylation specific factor 2) NDUFB1 (NADH dehydrogenase (ubiquinone) 1 beta subcomplex)	(GL397280.1: 29,388,208-29,433,598)	TRIP11 (thyroid hormone receptor interactor 11) FBLN 5 (fibulin 5)
<i>Otolemur gametii</i> <sup>a</sup>	SLC24A4 (Solute carrier family 24 (sodium/potassium/calcium exchanger), member 4) ENSOGAG00000032385 (cleavage and polyadenylation specific factor 2)	(GL873539.1: 18,068,653-18,100,321)	TRIP11 (thyroid hormone receptor interactor 11) FBLN 5 (fibulin 5)

**B - ATXN3L1**

Species	Flanking 5'	Sequence location	Flanking 3'
<i>Homo sapiens</i> <sup>a</sup>	TCEANC (transcription elongation factor A (SII) N-terminal and central domain containing) GPX1P1 (glutathione peroxidase pseudogene 1 provided) <sup>b</sup> EGFL6 (epidermal growth factor-like protein 6)	Xp.22.2	FAM9C (family with sequence similarity 9, member C) TMSB4X (thymosin beta 4, X-linked) TLR8 (Toll Receptor 9 Precursor, CD289 Antigen)
<i>Pan troglodytes</i> <sup>a</sup>	TCEANC (transcription elongation factor A (SII) N-terminal and central domain containing) EGFL6 (epidermal growth factor-like protein 6)	ChrX: 13,251,386-13,252,912	LOC735851 (family with sequence similarity 9, member C) LOC465493 (uncharacterized) B3Y655_PANTR (Toll Receptor 9 Precursor, CD289 Antigen)
<i>Gorilla gorilla</i> <sup>c</sup>	TCEANC (transcription elongation factor A (SII) N-terminal and central domain containing) EGFL6 (epidermal growth factor-like protein 6)	ChrX: 13240776-13241315	FAM9C (family with sequence similarity 9, member C) TMSB4X (thymosin beta 4, X-linked)
<i>Pongo abelii</i> <sup>a</sup>	TCEANC (transcription elongation factor A (SII) N-terminal and central domain containing) LOC100461104 (uncharacterized) LOC100460736 (epidermal growth factor-like protein 6)	ChrX: 13,230,000-13,231,067	LOC100452699 (protein FAM9C-like) TYB4_PONAB (thymosin beta 4, X-linked) TLR8 (Toll Receptor 9 Precursor, CD289 Antigen)
<i>Macaca mulatta</i> <sup>c</sup>	LOC711362 (transcription elongation factor A (SII) N-terminal and central domain containing) EGFL6 (epidermal growth factor-like protein 6)	ChrX: 11,016,157-11,017,194	LOC711039 (family with sequence similarity 9, member C) LOC710959 (Unknown) ENSMUG00000032533 (Unknown) ENSMUG00000032532 (Unknown) B6CK00_MACMU (Toll Receptor 9 Precursor, CD289 Antigen)
<i>Callithrix jacchus</i> <sup>a</sup>	LOC100388682 (transcription elongation factor A (SII) N-terminal and central domain containing) EGFL6 (epidermal growth factor-like protein 6)	ChrX: 11,295,476-13,523,142	LOC100415563 (thymosin beta 4, X-linked) TLR8 (Toll Receptor 9 Precursor, CD289 Antigen)
<i>Nomascus leucogenys</i> <sup>a</sup>	TCEANC (transcription elongation factor A (SII) N-terminal and central domain containing) EGFL6 (epidermal growth factor-like protein 6)	GL397281.1: 10,640,111-10,641,175	FAM9C (family with sequence similarity 9, member C) TLR8 (Toll Receptor 9 Precursor, CD289 Antigen)

**C - ATXN3L2**

Species	Flanking 5'	Sequence location	Flanking 3'
<i>Homo sapiens</i> <sup>b</sup>	KCNV1 (potassium channel, subfamily V, member 1) <sup>c</sup> RPSAP48 (ribosomal protein SA pseudogene 48)*	Chr8q23.2 (Chr8: 111567628-111568736)	LOC100129370 (NADH dehydrogenase (ubiquinone) 1 beta subcomplex) - pseudogene EEF1A1P37 (eukaryotic translation elongation factor 1 alpha 1 pseudogene 37)
<i>Pan troglodytes</i>	unknown		unknown
<i>Gorilla gorilla</i> <sup>c</sup>	KCNV1 (potassium channel, subfamily V, member 1) ENSGGOG00000022520 (40S Ribosomal)	Chr8: 109925669-109926225	ENSGGOG00000023880 (Unknown) ENSGGOG00000035637 (U4 spliceosomal RNA)* CSMD3 (CUB and Sushi multiple domains 3)*
<i>Pongo abelii</i> <sup>c</sup>	KCNV1_PONAB (potassium channel, subfamily V, member 1)*	Chr8: 117579303-117579866	CSMD3 (CUB and Sushi multiple domains 3)*
<i>Nomascus leucogenys</i> <sup>c</sup>	ENSNLEG00000024811 (U2 spliceosomal RNA)* KCNV1 (potassium channel, subfamily V, member 1)	GL397267:24598565-24599123	7SK (7SK RNA) ENSNLEG00000022484 (U4 spliceosomal RNA)* CSMD3 (CUB and Sushi multiple domains 3)*
<i>Macaca mulatta</i> <sup>c</sup>	ENSMUG00000034903 (U2 spliceosomal RNA) KCNV1 (potassium channel, subfamily V, member 1)	Chr8:112963801-112964361	LOC700556 (unknown) LOC699756 (Nucleophosmin) CSMD3 (CUB and Sushi multiple domains 3)*
<i>Callithrix jacchus</i> <sup>c</sup>	5S_rRNA (5S ribosomal RNA) LOC100389980 (Lysophospholipase 1)*	Chr19: 6218140-6218664	U3 (Small nucleolar RNA U3) C1orf143 (Unknown) TGFB2 (Transforming Growth Factor Beta 1 Precursor)

## 1.2. Phylogeny of *ATXN3* paralogues in primates

To assess the phylogenetic relationships of *ATXN3* paralogues identified in primates, we calculated the genetic distances among them (Table A1 of the appendix section). Based on these distances, a phylogenetic tree was constructed (Figure 8). Although we could not reproduce the phylogenetic tree of species for each one of the paralogues, we could observe that *ATXN3L1* and *ATXN3L2* orthologues clustered to each other, confirming their conservation along the primate lineage.

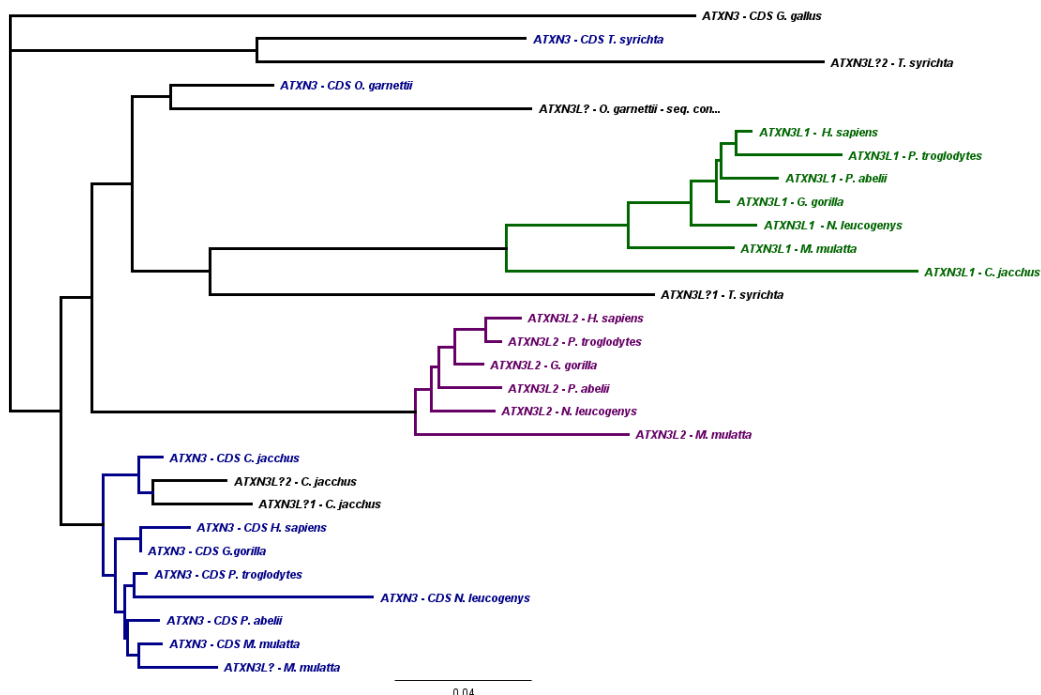


Figure 8 – Phylogenetic tree of *ATXN3*, *ATXN3L1* and *ATXN3L2* for primates, based on the genetic distances presented on Table A1 of the Appendix section. *Gallus gallus* was used as outgroup.

Marmoset *ATXN3L?2* was shown to be similar to its *ATXN3* coding sequence by sharing a recent node on the phylogenetic tree. This fact confirmed our reclassification of marmoset *ATXN3L2* into *ATXN3L?2* in the previous topic.

*ATXN3L?1* sequence of *Tarsius syrichta* was found to cluster next to *ATXN3L1* orthologues and to share more similarity with these orthologues than with the parental or *ATXN3L2* gene; for this reason, we reclassified it as *ATXN3L1*. On the other hand, the other Tarsier sequence, *ATXN3L?2*, shared a more recent ancestor with its parental sequence; therefore, its classification as *ATXN3L?* was maintained. As for *Otolemur garrattii*, *ATXN3L?* sequence did not seem to be orthologue of

*ATXN3L1* or *ATXN3L2*; instead, and as it displayed also less genetic distance to its parental gene, *ATXN3L?* classification was maintained.

The obtained phylogenetic tree was compared to others built based on *ATXN3*, *ATXN3L1* and *ATXN3L2* sequences (1) without the (CAG)<sub>n</sub> tract, and (2) containing the Josephin domain alone (Figures A2 and A3 of the Appendix section). No differences were found on the main clusters and nodes previously analysed to identify the origin of *ATXN3L1* and *ATXN3L2*.

After identifying *ATXN3L1* and *ATXN3L2* orthologues in primates, we wanted to gain insight into the events underlying their birth, namely the gene source. We, thus, calculated percentages of pairwise identities and identical sites by using aligned sequences of primates in which both retrocopies had previously been identified (Table 9).

**Table 9 - Pairwise identities and identical sites between *ATXN3* paralogues of several primate species. The coding sequence of the parental gene and the corresponding aligned sequences of *ATXN3L1* and *ATXN3L2* were used in these calculations. A single value is shown when both percentages are identical.**

	Pairwise identity (%)/identical sites (%)					
	<i>H. sapiens</i>	<i>P. troglodytes</i>	<i>G. gorilla</i>	<i>P. abelii</i>	<i>N. leucogenys</i>	<i>M. mulatta</i>
<i>ATXN3</i> vs <i>ATXN3L1</i>	77.8	78.2	75.6/79.2	77.7/79.6	76.5	80.0/80.5
<i>ATXN3</i> vs <i>ATXN3L2</i>	83.5	83.4	78.1/81.3	80.9/82.5	80.7	80.5/80.9
<i>ATXN3L1</i> vs <i>ATXN3L2</i>	71.1	69.8	74.1	70.7/70.8	73.0	71.2

A common pattern was observed for all species considered in this analysis: *ATXN3L2* presented a higher sequence identity to the *ATXN3* coding sequence than *ATXN3L1*. In addition, both *ATXN3L1* and *ATXN3L2* shared a higher sequence identity to the parental gene than between them, suggesting that two independent events of retrotransposition have occurred, instead of one event followed by duplication.

### 1.3. Identification of *ATXN3* retrocopies in other mammals

After clarifying the origin of *ATXN3* retrocopies in primates, we performed similar analyses for non-primate mammals to estimate the onset of the two most ancient retrotransposition events. Sequences homologous to *ATXN3L1* and *ATXN3L2* were obtained for *Pteropus vampyrus* (macrobat), *Cavia Porcellus* (guinea pig), *Oryctolagus cuniculus* (rabbit), *Canis familiaris* (dog), *Cloioepus hoffmanni* (sloth), *Bos taurus* (cow) and *Monodelphis domestica* (opossum). As none of these sequences was assigned as *ATXN3L1* or *ATXN3L2*, we named them as *ATXN3L?*. Table 10 resumes the main features of these sequences.

Table 10 – Additional *ATXN3* paralogues found in 7 non-primate mammals.

***ATXN3L?***

Species	Source	Location/Accession
<i>Pteropus vampyrus</i> (seq1)	BLAST - Trace Archives	gnl   ti:1371399947
<i>Pteropus vampyrus</i> (seq2)	BLAST - Trace Archives	gnl   ti:1328204199
<i>Cavia porcellus</i>	BLAST - Trace Archives	scaffold_33: 11,426,907-11,427,913
<i>Oryctolagus cuniculus</i>	BLAST - Trace Archives	gnl   ti:1979192258
<i>Canis familiaris</i> (seq1)	Ensembl	chr3: 10550724-10551819
<i>Canis familiaris</i> (seq2)	BLAST	NW_003726057.1 (chr5)
<i>Choloepus hoffmanni</i> (seq1)	BLAST - Trace Archives	gnl   ti:1361780914
<i>Choloepus hoffmanni</i> (seq2)	BLAST - Trace Archives	gnl   ti:1314261250
<i>Bos taurus</i> (seq1)	BLAST	chr8: 63086126-63086972
<i>Bos taurus</i> (seq2)	BLAST	chr10: 8893215-8893754
<i>Bos taurus</i> (seq3)	BLAST	chr4: 3842581-3843107
<i>Monodelphis domestica</i>	BLAST	chr5: 24707322-24707854

For the majority of the copies, here named *ATXN3L?*, size and location were not homogeneous among them. As for their synteny, we could not determine chromosomal location in cases where searches were done in the trace archives of NCBI since only scaffolds were available (*Pteropus vampyrus*, *Oryctolagus cuniculus* and *Choloepus hoffmanni*). For species annotated in the UCSC Genome Browser, however, sequences were localized through BLAT algorithm to obtain the chromosomal region and further analysed for the respective neighbouring genes. Therefore, in addition to the previous referred *ATXN3L?* of rhesus monkey and marmoset, synteny was analysed for dog, guinea pig, cow and opossum retrocopies (Table 11). None of these copies shared any flanking genes with *ATXN3*, *ATXN3L1* or *ATXN3L2*; instead, they differed also from each other. This data suggested that *ATXN3L?* sequences found in these mammalian species are in fact independent-origin retrocopies of *ATXN3*, and that none of the non-primate mammal sequences is orthologue of *ATXN3L1* or *ATXN3L2*.

**Table 11 - Synteny of ATXN3L? copies through evidence based on collinearity of genes in non-primate mammals. Two or three genes were retrieved within an interval of 600000 bp, except those marked with an asterisk (\*) which were found outside this interval. a, b and c next to the species name represent the source of the information: Ensembl, NCBI and UCSC Genome Browser, respectively. When information for a specific gene or location was retrieved from a source different from the indicated at species name, a, b or c specify the source.**

<i>ATXN3L?</i>				
Organism	Flanking 5'	Sequence location	Flanking 3'	Human homologues of flanking genes <sup>a</sup>
<i>Callithrix jacchus</i> <sup>a</sup>	<i>FA2H</i> (fatty acid 2-hydroxylase) <i>MLKL</i> (mixed lineage kinase domain-like)	Chr20: 28,630,500-28,631,553	<i>WDR59</i> (WD Repeat containing 59) <i>ENSCJAG00000014388</i> (Class I antigen)	Chr16
<i>Macaca mulata</i> <sup>b</sup>	<i>P2RX7</i> (P2X Purinoceptor 7) <i>ENSMUG00000031514</i> (Unknown)	Chr11: 122178105-122178311	<i>SNORA70</i> (Small nucleolar RNA <i>SNORA70</i> ) <i>OASL</i> (2'-5'-oligoadenylate synthetase-like) <i>LOC701243</i> (Hepatocyte nuclear factor 1 alpha)	Chr12
<i>Canis familiaris</i> <sup>a</sup>	<i>NUDT12</i> (Peroxisomal NADH pyrophosphatase) <i>GIN1</i> (Gypsy Retrotransposon Integrase 1) <i>ENSCAFG00000024608</i> (Unknown)	Chr3: 10550724-10551819	<i>ENSCAFG00000002024</i> (Heterogeneous Nuclear Ribonucleoprotein) <i>CDH1</i> (chromodomain helicase DNA binding protein 1)	Chr5
<i>Cavia porcellus</i> <sup>c</sup>	<i>ENSCPOG00000026477</i> (miRNA) <i>ENSCPOG00000016622</i> (5S_rRNA)	scaffold_33: 11,426,907-11,427,913	<i>ENSCPOG000000021308</i> (Scavenger Receptor Cystein Rich Type 1 M130 Precursor) <i>ENSCPOG00000009339</i> (Unknown)*	Unknown
Organism	Flank	Sequence location	Flank	Human homologues of flanking genes <sup>a</sup>
<i>Canis familiaris</i> <sup>c</sup>	<i>ELAC2</i> (ElaC homolog 2 - Zinc phosphodiesterase)* <i>5S_rRNA</i> (5S ribosomal RNA) <i>ENSCAFG00000017903</i> (Heparan sulfate glucosamine 3-O-sulfotransferase)	Chr5: 40677963-40679029	<i>COX10</i> ( <i>COX10</i> homolog, cytochrome c oxidase assembly protein) <i>ENSCAFG00000017893</i> (Heparan sulfate glucosamine 3-O-sulfotransferase)	Chr17
<i>Bos taurus</i> <sup>c</sup>	<i>C9orf174</i> - Chromosome 9 open reading frame 174 (uncharacterized) hypothetical protein <i>LOC507550</i> (uncharacterized)	Chr8: 63086126-63086972	<i>TDRD7_BOVIN</i> (Tudor domain-containing protein 7) <i>TMOD1_BOVIN</i> (Tropomodulin-1)	Chr9
<i>Bos taurus</i> <sup>c</sup>	<i>SNORA31</i> (small nucleolar RNA, H/ACA box 31)	Chr4: 3842581-3843107	<i>A0JNG1</i> (Cordon Bleu Gene) <i>Q6U8D5_Bovin</i> (Growth Factor Receptor Bound Adapter)	Chr7
<i>Bos taurus</i> <sup>c</sup>	<i>OTP</i> (Orthopedia Homeobox) <i>TBCA</i> (Tubulin-specific chaperone A)	Chr10: 8893215-8893754	<i>AP3B1</i> (AP-3 complex subunit beta-1) <i>Q3T0D2_Bovin</i> (Secretory Carrier-Associated Membrane Protein 1)	Chr5
<i>Monodelphis domestica</i> <sup>c</sup>	<i>C4orf37</i> (Chromosome 4 open reading frame 379) <i>PPM1K</i> (protein Phosphatase, Mg2+/Mn2+ dependent, 1K)*	Chr5:24707322-24707854	<i>XM_001376187.1</i> (Casein Kinase I)*	Chr4

In order to confirm the origin of *ATXN3L?* copies, sequences were aligned with the coding regions of the respective parental gene, as well as with paralogue sequences previously identified for primates. From this overall alignment (Figure A1 of Appendix section), we calculated genetic distances (Table A4 of the Appendix section) and constructed a phylogenetic tree (Figure 9). Results showed all *ATXN3L?* sequences more similar to the respective parental genes than to primate *ATXN3L1* or *ATXN3L2* sequences. Indeed, they seemed to have originated independently and more recently after species divergence, instead of sharing common ancestral origins. In addition, the extra paralogue sequences found for *M. mulatta* (*ATXN3L?*) and *C. jacchus* (*ATXN3L?1* and *ATXN3L?2*) also showed a most common recent ancestor with their parental sequences than with the other paralogues. This suggests that all sequences classified as *ATXN3L?* resulted, most likely, from additional retrotransposition events originated, more recently, from the respective parental genes.

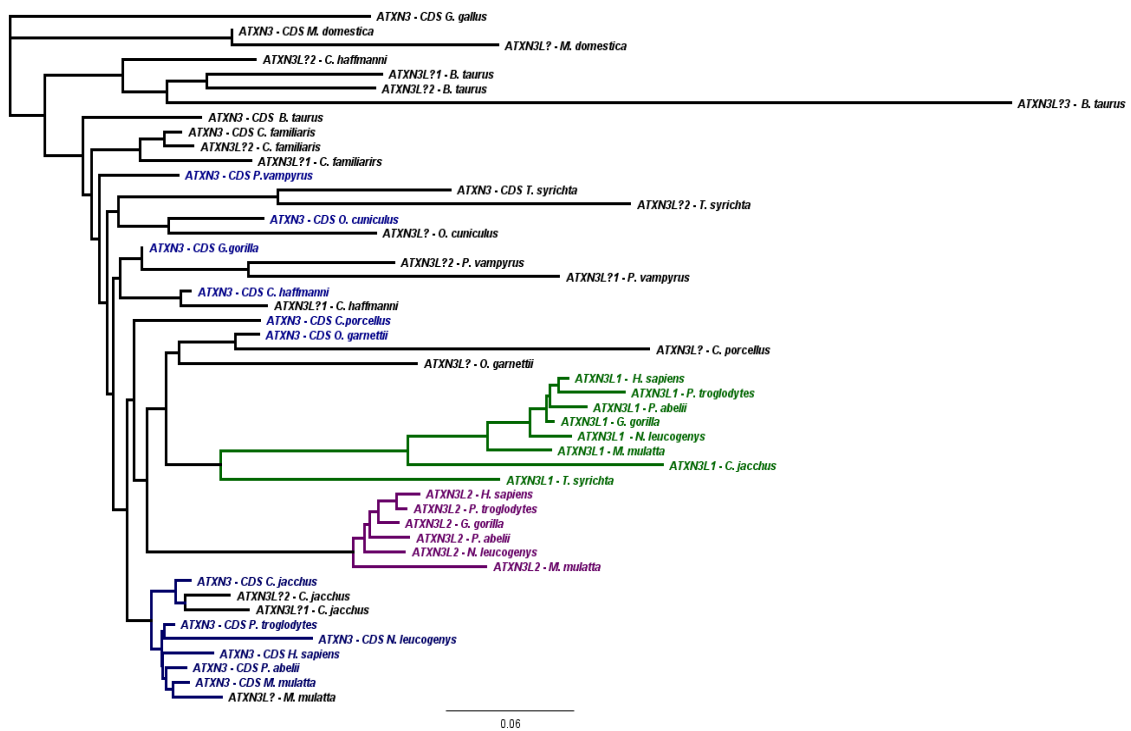
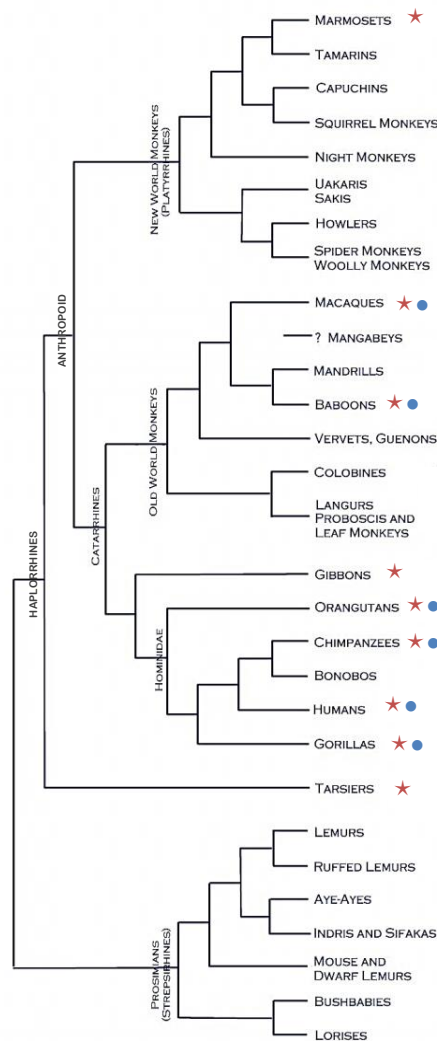


Figure 9 - Phylogenetic tree for *ATXN3*, *ATXN3L1*, *ATXN3L2* and *ATXN3L?* retrocopies of mammals, based on the genetic distances represented on table A4 of the Appendix section. *Gallus gallus* was used as outgroup.

To confirm these results, we constructed phylogenetic trees of these three genes for all mammalian species (1) without the (CAG)<sub>n</sub> tract (usually a source of much variation), and (2) for the high conserved Josephin domain alone (Figures A4 and A5). For the first, a similar tree was obtained, whereas the phylogenetic tree based exclusively on the sequence that encodes the Josephin domain (in *ATXN3* and homologous sequences in its paralogues), has shown *ATXN3L1* genetically closer to the parental *ATXN3* than *ATXN3L2*, suggesting that both have very similar Josephin domains, with most of the variation accumulated in the unstructured C-terminal.

Finally, after performing all previous studies, we were able to estimate the time for the origin of *ATXN3* paralogues (Figure 10), based on the tree of species and respective divergence times, in Figure A6 of the appendix section.



**Figure 10 – Phylogenetic tree of primates. The presence of *ATXN3L1* and *ATXN3L2* is marked in each species with a red star and a blue circle, respectively. For *ATXN3L1*, the absence of the red star in Haplorrhines may be explained by incomplete gene annotation in databases. The same applies to blue circles marking *ATXN3L2* in Catarrhini (adapted from <http://whozoo.org/mammals/Primates/primatephylogeny.htm>).**



*ATXN3L1* is present in the entire clade of Haplorrhines but absent in Strepsirrhines, indicating that the retrotransposition event on its origin occurred about 63 million years ago (MYA). As for *ATXN3L2*, since orthologues are present in the entire clade of Catarrhines but absent in Platyrrhines, the retrotransposition event is likely to have occurred about 35 MY ago.<sup>79</sup>

## **2. *ATXN3* transcripts involved in the retrotransposition events**

The alignment of *ATXN3L1* and *ATXN3L2* with the 21 human protein coding transcripts annotated in Ensembl has shown both paralogues more similar to *ATXN3-001* than to any other transcript. Based on these results, both events of retrotransposition on the origin of *ATXN3L1* and *ATXN3L2* seemed to have involved this transcript, since they all share exons 1, 3 (part), 4 (part), 6 (part), 7, 9, 11, 12, 13 (part), 20 (part) and 21 (based on Ensembl exon classification).

## **3. Selective signatures underlying *ATXN3L1* and *ATXN3L2***

Evolutionary pressures on proteins were quantified by calculating the omega ratio, or the ratio of substitution rates at non-synonymous and synonymous sites, for *ATXN3* and *ATXN3L1* genes, for different primate species (Table 12).

Table 12 – dN/dS (omega) ratio calculations for *ATXN3* and *ATXN3L1*, to illustrate the selective signatures underlying their evolution.

vs.		<i>ATXN3</i>			<i>ATXN3L1</i>		
		dN	dS	$\omega$	dN	dS	$\omega$
Hsap	Ptro	0,0026	0,0094	0,2766	0,0291	0,0549	0,5301
Hsap	Ggor	0,1146	0,1613	0,7105	0,0100	0,0246	0,4065
Hsap	Pabe	0,0026	0,0190	0,1368	0,0163	0,0601	0,2712
Hsap	Nleu	0,0051	0,0287	0,1777	0,0252	0,0979	0,2574
Hsap	Mmul	0,0051	0,0686	0,0743	0,0323	0,2266	0,1425
Hsap	Cjac	0,0129	0,1004	0,1285	0,1431	0,3902	0,3667
Hsap	Tsyr	0,6000	0,8995	0,6670	0,2099	0,5639	0,3722
Ptro	Ggor	0,1175	0,1500	0,7833	0,0304	0,0396	0,7677
Ptro	Pabe	0,0051	0,0094	0,5426	0,0389	0,0730	0,5329
Ptro	Nleu	0,0077	0,0191	0,4031	0,0421	0,1142	0,3687
Ptro	Mmul	0,0077	0,0585	0,1316	0,0507	0,2388	0,2123
Ptro	Cjac	0,0103	0,0900	0,1144	0,1605	0,4522	0,3549
Ptro	Tsyr	0,6270	0,8730	0,7182	0,2335	0,5919	0,3945
Ggor	Pabe	0,1146	0,1613	0,7105	0,0163	0,0345	0,4725
Ggor	Nleu	0,1160	0,1678	0,6913	0,0233	0,6830	0,0341
Ggor	Mmul	0,1190	0,1792	0,6641	0,0278	0,1919	0,1449
Ggor	Cjac	0,1174	0,2500	0,4696	0,1378	0,3792	0,3634
Ggor	Tsyr	0,1414	1,1120	0,1272	0,2022	0,5319	0,3801
Pabe	Nleu	0,0026	0,0287	0,0906	0,0291	0,0867	0,3356
Pabe	Mmul	0,0026	0,0483	0,0538	0,0375	0,1940	0,1933
Pabe	Cjac	0,0103	0,1004	0,1026	0,1501	0,3617	0,4150
Pabe	Tsyr	0,0572	0,8694	0,0658	0,2183	0,5221	0,4181
Nleu	Mmul	0,0051	0,0792	0,0644	0,0355	0,1980	0,1793
Nleu	Cjac	0,0129	0,0900	0,1433	0,1398	0,3840	0,3641
Nleu	Tsyr	0,0572	0,8745	0,0654	0,2126	0,5783	0,3676
Mmul	Cjac	0,0129	0,1561	0,0826	0,1334	0,4354	0,3064
Mmul	Tsyr	0,0614	0,7724	0,0795	0,1979	0,7145	0,2770
Cjac	Tsyr	0,0627	0,9714	0,0645	0,2529	0,7266	0,3481

Omega values were always less than unity, for both *ATXN3* and *ATXN3L1* in all species comparisons. However, the rate of variance of these values were higher in the parental gene (from approximately 0.06 to 0.80) than in the *ATXN3L1* paralogue, which showed less different values (mainly varying from 0.10 to 0.50). The fact that parental genes presented higher variance can be partly explained by the sequence differences found in gorilla, orangutan and tarsier. Still, as omega is less than unity in all calculations, our results suggested that selective constrains are suppressing protein changes, indicating that, as *ATXN3*, also the function of *ATXN3L1* has been kept in check.

#### 4. Nucleotide diversity of *ATXN3L1* and *ATXN3L2*

We have assessed genetic variation of all *ATXN3* paralogues by sequencing a total of 49 human DNA samples (Table 13). For *ATXN3*, we have found a total of seventeen variations, mostly in intronic regions (or in regions that can be transcribed only by less common alternative splicing), according to the Ensemble exon/intron sequence annotations. This is consistent with the fact that transcriptionally active regions of the genome are more prone to be conserved. In exonic regions, we have found six variations: five non-synonymous and one synonymous coding. From these six variations, however, only two (one synonymous and one non-synonymous) were common to all transcripts; the other four were present only in three or fewer transcripts, none of them on the four most common *ATXN3* transcripts.

**Table 13 – *ATXN3*, *ATXN3L1* and *ATXN3L2* nucleotide diversities. Base positions of *ATXN3* were determined based on the sequence annotation of Ensembl; for *ATXN3L1* and *ATXN3L2*, the A nucleotide of the initiation codon (the corresponding first AAG codon in the case of *ATXN3L2*) is base number 1. a – in exonic regions of *ATXN3*-206, 207 and 209 transcripts; b - in exonic regions of *ATXN3*-015 transcript; \*ancestral alleles determined from primate orthologues alignments; Amino acids: V - Valine; T – Threonine; E - Glutamic acid; R – Arginine; M- Methionine; K – Lysine; Y – Tyrosine; D - Aspartic acid; G – Glycine; NA – Not applicable.**

	Variation ID	Base position	Number of chromosomes analysed	Ancestral allele	Base change	Absolute frequencies	Relative frequencies	Annotated frequencies	Amino acid change
<b><i>ATXN3</i> (Chr. 14)</b>	rs17847278	160 (Intron 4)	20	G	A/G	7/13	0.35/0.65	0.29/0.71	NA
	rs1997920	194 (Exon 6)	48	T	T/C	24/24	0.50/0.50	0.31/0.69	NA
	rs1997919	340 (Exon 6)	42	G	G/A	22/20	0.52/0.48	0.69/0.31	NA
	rs1997918	417 (Exon 6)	42	A	G/A	3/39	0.07/0.93	0.01/0.99	NA
	rs1997917	434 (Exon 6)	44	G	G/A	23/21	0.48/0.52	0.27/0.73	NA
	rs4904834	524 (Intron6)	42	C*	G/C	24/18	0.57/0.53	0.45/0.55	NA
	rs8003520	1081 (Intron 8)	16	G	G/A	3/13	0.19/0.81	0.30/0.70	NA
	rs12590497	129 (Intron 10)	16	T	G/T	10/6	0.63/0.37	0.68/0.32	NA
	rs16999141	17 (Exon 11)	12	T	C/T	6/6	0.50/0.50	0.47/0.53	V/V
	rs1048755	26 (Exon 12)	16	G	A/G	6/10	0.31/0.69	0.26/0.74	M/V
	rs761553	1220 (Intron 12)	16	C	G/C	6/10	0.38/0.62	0.68/0.32	NA
	rs761552	31 (Exon 14)	14	C	T/C	7/7	0.50/0.50	0.32/0.68	M/T <sup>a</sup>
	rs761551	24 (Intron14)	14	C	A/C	7/7	0.50/0.50	0.31/0.69	NA
	rs10467858	127 (Intron 17)	12	G*	A/G	6/6	0.50/0.50	0.27/0.73	NA
	rs7158733	201 (Exon20)	18	C*	C/A	13/5	0.72/0.28	0.68/0.32	Y/Stop <sup>b</sup>
	rs3092822	261 (Exon 20)	22	A*	A/C	14/8	0.64/0.36	0.68/0.32	E/D <sup>b</sup>
	c.2926A>G	319 (Exon20)	22	A*	A/G	8/14	0.36/0.64	undetermined	K/E <sup>b</sup>
<b><i>ATXN3L1</i> (Chr. X)</b>	c.939G>A	939	38	G*	A/G	2/36	0.05/0.95	undetermined	R/R
c.995A>G	995	37	A*	A/G	6/31	0.16/0.84	undetermined	D/G	
<b><i>ATXN3L2</i> (Chr. 8)</b>	c.230T>C	230	68	T*	C/T	2/66	0.03/0.97	undetermined	NA
c.398G>A	398	76	G*	A/G	1/74	0.01/0.99	undetermined	NA	
c.752G>A	752	64	G*	A/G	3/61	0.04/0.96	undetermined	NA	

For *ATXN3L1*, only two variations were detected, both near the end of the putative coding sequence: one synonymous and the other non-synonymous, but both relatively rare, varying from 0.05 to 0.16%, respectively. Considering *ATXN3L2*, three variations were observed along the sequence; however, these variants were rare, among our analysed samples, with percentages lower than 0.04%. We have notice, however, that not all SNPs annotated in databases were detected by us, in the resulting sequences, probably due to the low number of chromosomes analysed.

## 5. Evolution of the (CAG)<sub>n</sub> tract in *ATXN3*, *ATXN3L1* and *ATXN3L2*

To gain insight into the processes by which the (CAG)<sub>n</sub> tracts have been accumulating variation throughout each paralogue evolution, we analysed sequences from all primates obtained from databases (Table 14). However, this analysis involved mainly a single sequence from each species, which did not allow us to study *loci* diversity. To analyse the level of polymorphism, *ATXN3L1* and *ATXN3L2* from 6 primate species were sequenced (Table 15) and data gathered with the reference ones.

Table 14 – *ATXN3*, *ATXN3L1* and *AtxN3L2* (CAG)<sub>n</sub> tracts of several primates collected from NCBI, Ensembl and UCSC Genome Browser databases.

	<i>ATXN3</i>	<i>ATXN3L1</i>	<i>ATXN3L2</i>
<i>Homo sapiens</i>	(CAG) <sub>2</sub> CAA AAG CAG CAA (CAG) <sub>17</sub> Q <sub>2</sub> Q K Q Q Q <sub>17</sub>	(CAG) <sub>2</sub> GAA CAG AAG (CAG) <sub>2</sub> (CAA) <sub>2</sub> CAG Q <sub>2</sub> E Q K Q <sub>2</sub> Q <sub>2</sub> Q	(CAG) <sub>3</sub> (CGG CAG) <sub>n</sub>
<i>Pan troglodytes</i>	CAG CAA (CAG) <sub>12</sub> Q Q Q <sub>12</sub>	(CAG) <sub>2</sub> GAA CAG AAG (CAG) <sub>2</sub> CAA CAC CAG Q <sub>2</sub> E Q K Q <sub>2</sub> Q H Q	CAG (CGG CAG) <sub>9</sub>
<i>Gorilla gorilla</i>	(CAG) <sub>2</sub> CAA AAG CAG CAA AAG (CAG) <sub>n</sub> Q <sub>2</sub> Q K Q Q K Q <sub>n</sub>	(CAG) <sub>2</sub> GAA CAG AAG (CAG) <sub>2</sub> (CAA) <sub>2</sub> CAG Q <sub>2</sub> E Q K Q <sub>2</sub> Q <sub>2</sub> Q	(CAG) <sub>10</sub>
<i>Pongo abelii</i>	(CAG) <sub>2</sub> CAA AAG CAG CAA (CAG) <sub>8</sub> Q <sub>2</sub> Q K Q Q Q <sub>8</sub>	(CAG) <sub>2</sub> GAA CAG AAG CTG (CAG) <sub>4</sub> Q <sub>2</sub> E Q K L Q <sub>4</sub>	(CAG) <sub>6</sub>
<i>Nomascus leucogenys</i>	(CAG) <sub>3</sub> CAA (CAG) <sub>9</sub> Q <sub>3</sub> Q Q <sub>9</sub>	(CAG) <sub>2</sub> GAA CAG AAG (CAG) <sub>2</sub> (CAA) <sub>2</sub> CAG Q <sub>2</sub> E Q K Q <sub>2</sub> Q <sub>2</sub> Q	(CAG) <sub>8</sub>
<i>Macaca mulatta</i>	(CAG) <sub>2</sub> CAA (CAG) <sub>2</sub> AAG (CAG) <sub>7</sub> Q <sub>2</sub> Q Q <sub>2</sub> K Q <sub>7</sub>	(CAG) <sub>2</sub> GAA CAG AAG (CAG) <sub>2</sub> (CAA) <sub>2</sub> CAG Q <sub>2</sub> E Q K Q <sub>2</sub> Q <sub>2</sub> Q	(CAG) <sub>6</sub> CAC CAG
<i>Callithrix jacchus</i>	CAG CAA (CAG) <sub>2</sub> CAA (CAG) <sub>2</sub> [GAG (CAG) <sub>3</sub> ] <sub>2</sub> Q Q Q <sub>5</sub> Q Q <sub>n</sub> (E Q <sub>3</sub> ) <sub>2</sub>	CAA GAA CAA (CAG) <sub>2</sub> Q E Q Q <sub>2</sub>	Not described
<i>Tarsius syrichta</i>	Not completely sequenced	CAG CAA (CAG) <sub>2</sub> Q Q Q <sub>2</sub>	Not described
<i>Otolemur garnettii</i>	(CAG) <sub>5</sub> CAA (CAG) <sub>3</sub> AGG Q <sub>5</sub> Q Q <sub>3</sub> K	Not described	Not described

Analysing the compiled sequences alone, (CAG)<sub>n</sub> tract was found to be more conserved in the parental gene along the primate lineage, mainly in human, gorilla and orangutan, with a (CAG)<sub>2</sub> CAA AAG CAG CAA tract, followed by AAG and/or (CAG)<sub>n</sub>. This configuration gives rise to an almost pure polyQ tract, with a single K (lysine) interruption (two in gorilla). Chimp, on the other hand, presents a more simple CAG CAA (CAG)<sub>n</sub> configuration, encoding a pure glutamine stretch. The (CAG)<sub>n</sub> tract in gibbon, rhesus monkey, marmoset and bushbaby has a (CAG)<sub>n</sub> CAA (CAG)<sub>n</sub>, followed by a variable end. The protein tract is pure (in gibbon) or interrupted only by a lysine in rhesus monkey and bushbaby, or a glutamic acid in marmoset. For *ATXN3L1*, the (CAG)<sub>n</sub> is also conserved among primates, having a (CAG)<sub>2</sub> GAA CAG AAG (CAG)<sub>2</sub> stretch followed by a small variable end with CAA, CAG or CAC triplets. The resulting putative protein sequence is interrupted, generally with the Q<sub>2</sub> E Q K Q<sub>n</sub> configuration. Marmoset and tarsier have a different configuration, with a smaller tract: CAA GAA CAA (CAG)<sub>2</sub> and CAG CAA (CAG)<sub>2</sub>, respectively. If translated, these tracts give rise to a smaller pure (for tarsier) or almost pure (for marmoset)

polyglutamine stretches, which may suggest that ATXN3L1 CAG repeat has expanded and gained interruptions along the primate lineage. As for ATXN3L2, the (CAG)<sub>n</sub> tract is more variable. For gorilla, orangutan and gibbon, the (CAG)<sub>n</sub> tract is pure, whereas for human and chimp it is highly interrupted by CGG codons. For the first time, a hexanucleotide repeat is observed associated within ataxin-3 paralogues, instead of a trinucleotide pattern. These results suggested that ATXN3L1 acquired interruptions in its CAG repeat region that may turn it more stable. On the other hand, in the case of ATXN3L2, the almost pure (CAG)<sub>n</sub> in species which diverged earlier in primate evolution, and the polymorphic (CGG CAG)<sub>n</sub> in humans and chimps has shown the highest instability associated to this locus.

**Table 15 - ATXN3, ATXN3L1 and ATXN3L2 (CAG)<sub>n</sub> tracts obtained by sequencing for several primates.\* mark the minimum number of chromosomes analysed taking into account that we ignored whether the respective gene location in these species were in autosomal or sex chromosomes (which, in case of homozygous regions could result in different number of alleles analysed: n vs 2n.**

	ATXN3L1			ATXN3L2		
	Chr.	Number of chromosomes	Sequence	Chr.	Number of chromosomes	Sequence
<i>Homo sapiens</i>	X	36	(CAG) <sub>2</sub> GAA CAG AAG (CAG) <sub>2</sub> (CAA) <sub>2</sub> CAG Q <sub>2</sub> E Q K Q <sub>2</sub> Q <sub>2</sub> Q	8	1 (CAG) <sub>3</sub> (CGG CAG) <sub>11</sub> 4 (CAG) <sub>3</sub> (CGG CAG) <sub>9</sub> 38 (CAG) <sub>3</sub> (CGG CAG) <sub>8</sub> 3 (CAG) <sub>3</sub> (CGG CAG) <sub>6</sub> 2 (CAG) <sub>3</sub> (CGG CAG) <sub>5</sub> 10 (CAG) <sub>3</sub> (CGG CAG) <sub>4</sub>	
<i>Pan troglodytes</i>	X	3	(CAG) <sub>2</sub> GAA CAG AAG CTG CAG (CAA) <sub>2</sub> CAG Q <sub>2</sub> E Q K L Q Q <sub>2</sub> Q	8	3 CAG (CGG CAG) <sub>9</sub> 1 CAG (CGG CAG) <sub>6</sub>	
<i>Gorilla gorilla</i>	X	2	(CAG) <sub>2</sub> GAA CAG AAG CTG CAG (CAA) <sub>2</sub> CAG Q <sub>2</sub> E Q K L Q Q <sub>2</sub> Q	8	1 (CAG) <sub>16</sub> 1 (CAG) <sub>10</sub> 2 (CAG) <sub>7</sub> 2 (CAG) <sub>6</sub>	
<i>Pongo abelii</i>	X	2	(CAG) <sub>2</sub> GAA CAG AAG CTG CAG (CAA) <sub>2</sub> CAG Q <sub>2</sub> E Q K L Q Q <sub>2</sub> Q	8	1 (CAG) <sub>7</sub> 1 (CAG) <sub>6</sub>	
<i>Papio</i>	?			?	1* (CAG) <sub>6</sub> CAC CAG 2* (CAG) <sub>7</sub> CAC CAG 1* CAG CAT (GAG) <sub>5</sub> CAC CAG	
<i>Macaca mulatta</i>	X	2	(CAG) <sub>2</sub> GAA CAG AAG (CAG) <sub>2</sub> (CAA) <sub>2</sub> CAG Q <sub>2</sub> E Q K Q <sub>2</sub> Q <sub>2</sub> Q	8	4 (CAG) <sub>6</sub> CAC CAG 4 (CAG) <sub>6</sub> CAC CAG	
<i>Macaca fascicularis</i>	?	2*	(CAG) <sub>2</sub> GAA CAG AAG (CAG) <sub>2</sub> (CAA) <sub>2</sub> CAG Q <sub>2</sub> E Q K Q <sub>2</sub> Q <sub>2</sub> Q	?	3* (CAG) <sub>6</sub> CAC CAG	

Comparing the compiled sequences with those obtained by us through sequencing (Table 15), we noticed that ATXN3L1 sequences of some species may not be correctly described in databases. For chimp, gorilla and orangutan, observed CAG repeat configurations differed from those described in Table 14. This could be simply explained by the level of polymorphism associated to these tracts; nevertheless, all repeat tracts sequenced by us for individuals from these 3 species presented no variation, all sharing the same configuration and length: (CAG)<sub>2</sub> GAA CAG AAG CTG

CAG (CAA)<sub>2</sub> CAG, possibly encoding a Q<sub>2</sub> E Q K L Q<sub>4</sub> protein stretch. This means that this region may be even more conserved along the primate lineage than it has been described. For *ATXN3L2*, the (CAG)<sub>n</sub> tract of all sequenced primate samples showed a similar configuration to reference sequences available in databases, illustrated in Table 14; the repeat size showed, however, to be more variable.

## **Part 2 – Insights into *ATXN3L1* and *ATXN3L2* functional relevance**

### **1. ORF predictions for *ATXN3L1* and *ATXN3L2***

To explore the possibility of *ATXN3L1* and *ATXN3L2* paralogues being transcribed, the analysis of the respective ORFs was performed; the obtained results are illustrated in Figure 11. *ATXN3L1* showed a complete and intact ORF that is maintained in the 7 primate species previously analysed by us. These results confirmed the hypothesis that this paralogue is likely to be subjected to selective constraints that are keeping its sequence on check and almost intact along the primate lineage. On the other hand, for *ATXN3L2*, there are no intact ORFs predicted in any of the primate sequences; instead, only short ORFs of approximately 200 bp or less are expected.

**Evolution and Functional Relevance of Ataxin-3 Paralogues**  
 Maria Inês Martins

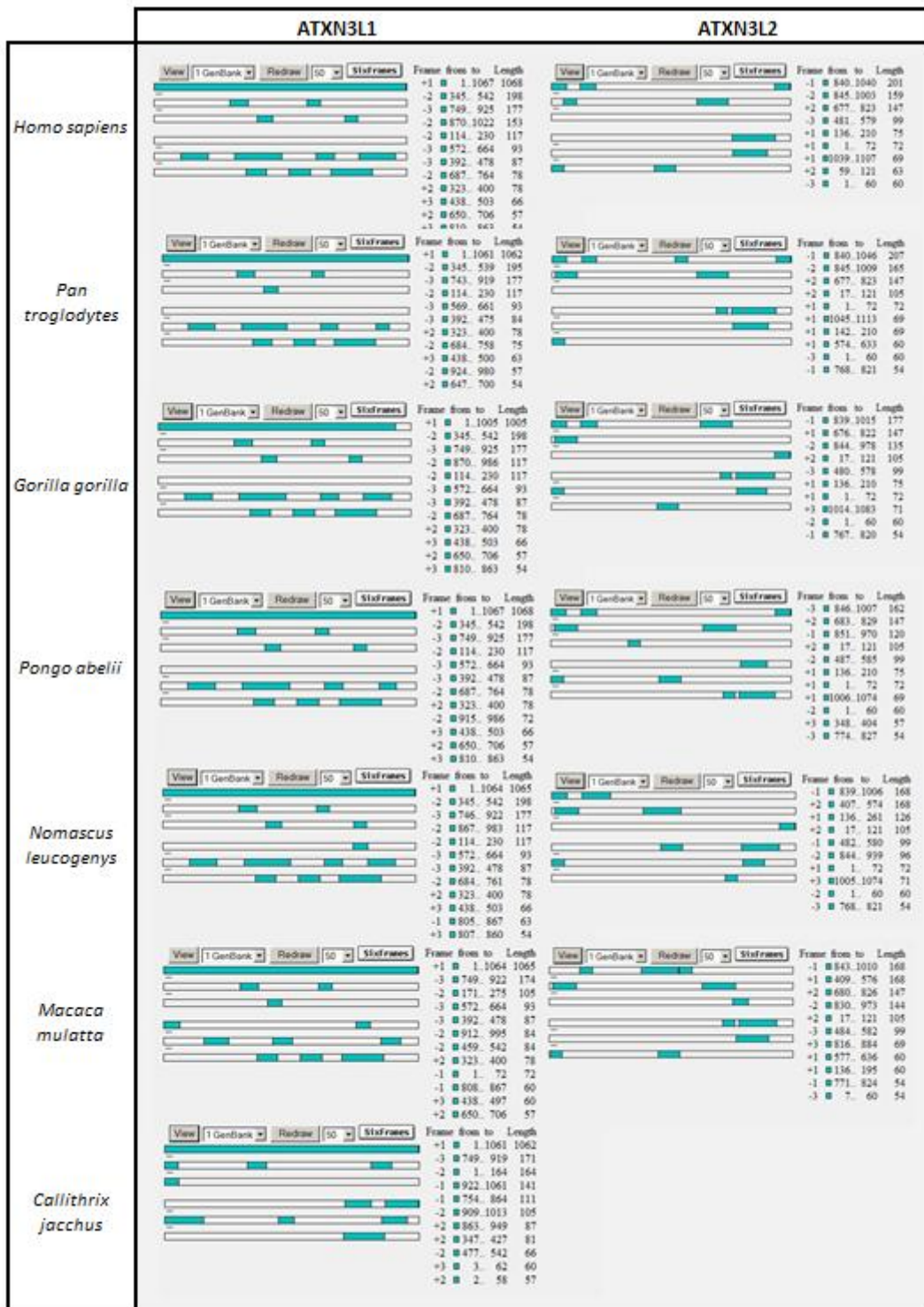


Figure 11 - Predicted ORFs for *ATXN3L1* and *ATXN3L2* primate sequences, calculated by NCBI ORF finder. Obtained ORF predictions are illustrated in blue and the corresponding sizes are described on the right side of each illustration.

## 2. Analysis of ATXN3L1 putative protein domains in comparison to the parental ATXN3

To understand how differences on the *ATXN3L1* sequence relatively to the parental gene could possibly alter its putative coding protein, the principal domains and regions important to ataxin-3 functions were compared between the two sequences (Figure 12).

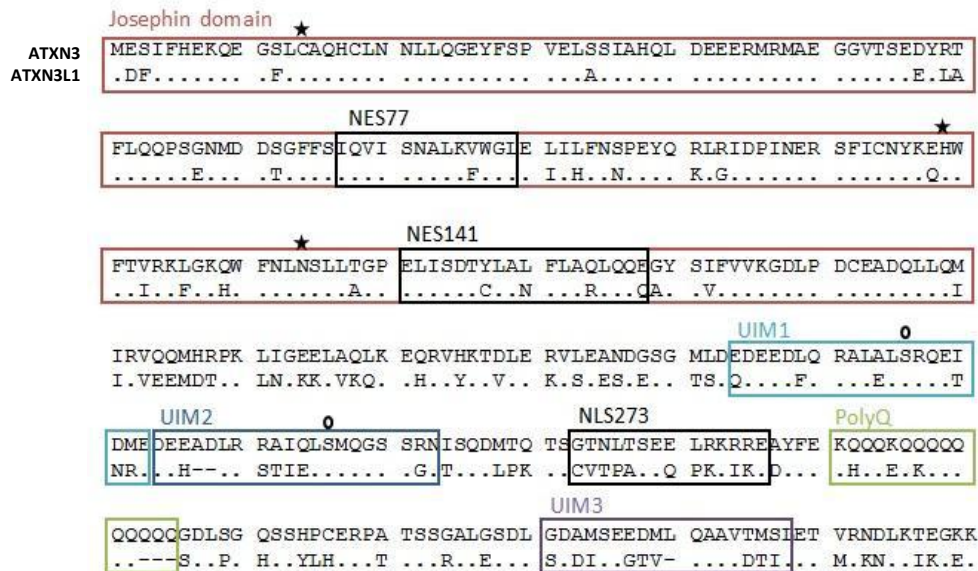


Figure 12 – Comparison between ATXN3 and ATXN3L1 protein coding sequences. Dots show conserved amino acids and letters show the amino acid changes. The main functional Josephin domain, the ubiquitin interacting motifs (UIMs) and polyQ tract are highlighted. Other important sites of the protein are also represented, as the nuclear export signal (NES) and the nuclear localization signal (NLS). ★ marks amino acids composing the catalytic triad; ○ marks the conserved serine-232 and serine-260 in UIM1 and UIM2, respectively.

The catalytic Josephin domain was found to be highly conserved in ATXN3L1, with the important catalytic triad, composed of C14, H119 and N134, totally preserved. In addition, NES77 and NES141 are composed mostly by the same amino acids. Additionally, UIM1 of ATXN3L1 is quite similar to the parental protein, although the same cannot be observed for UIM2 and UIM3, which possess more amino acid changes and some gaps. The NLS273 is highly different when compared to the parental protein. Finally, the polyQ tract, as described above, has some alterations that turned it more interrupted in ATXN3L1. It is important to understand that despite these two protein sequences present some amino acid divergences, these small alterations may not correspond to a significant change on the protein conformation and/or function. To better understand this divergence features will be necessary to proceed with more functional studies to characterize the putative ATXN3L1 protein sequence.



### 3. mRNA expression profile of *ATXN3L1* in humans

We tested for the presence of *ATXN3L1* in cDNA from different human body tissues by using primers designed to specifically amplify this putative transcript. First, to overview the transcriptional pattern among the 16 selected tissues, two of the four planned amplifications were done at a lower temperature (56°C). Primers for the amplification of actin cDNA were used as control (Figure 13).

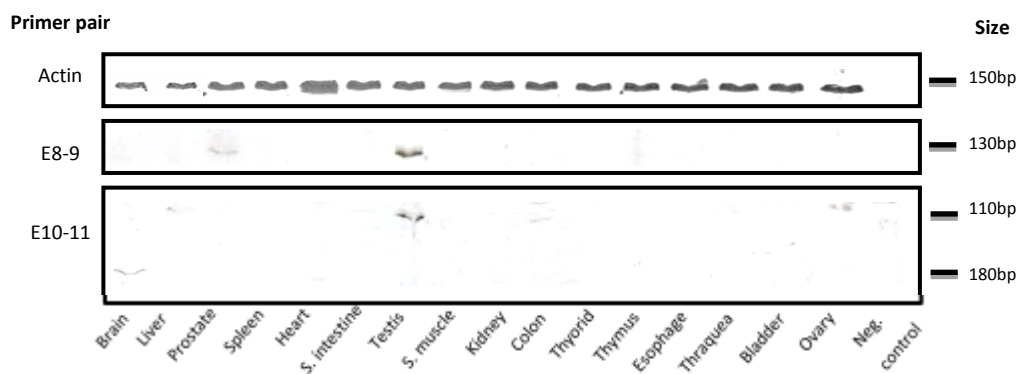


Figure 13 – *ATXN3L1* transcriptional pattern of 16 human tissues using E8-9 and E10-11 primers. The human tissues used are listed at the bottom of the image and the sizes of obtained bands are labeled on the right side. Actin cDNA was amplified as a control.

With these conditions, we have found discordant results for the two amplifications tested. By using E8-9 primers, amplification was obtained in only two tissues (prostate and testis), whereas for E10-11, amplified products were observed in brain, liver, testis, colon and ovary. In both cases, testis cDNA was the tissue for which we observed more amplified product, allowing the formation of a clearly visible gel band (a similar amplification pattern for actin cDNA in all tissues discarded the possible heterogeneity among different tissue samples). Sequencing of each band showed that all amplified bands were specific human *ATXN3L1* sequences. Additionally, E10-11 primers amplified a longer fragment (180 bp) than the expected 110 bp on brain cDNA; we still need to clarify the reason for this difference on fragment sizes.

To increase the specificity of *ATXN3L1* cDNA amplification, the experiment was repeated by applying new different condition of amplification. We increased the annealing temperature to 58°C on the PCR reaction and tested the amplification with the four designed pairs of primers, on the 8 cDNA samples for which we had previously observed amplification (Figure 13).

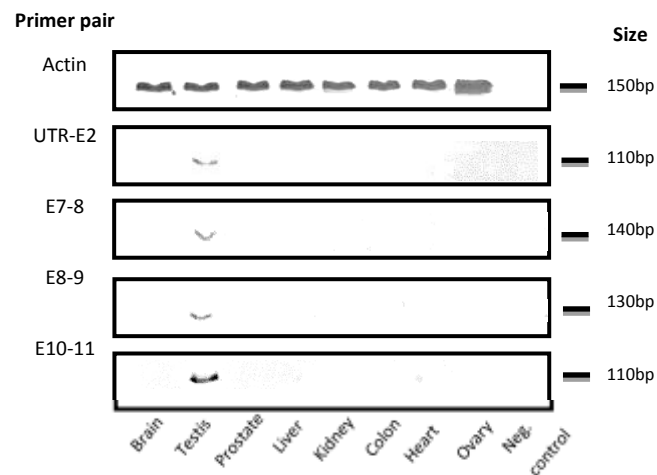


Figure 14 - ATXN3L1 transcriptional pattern for the 4 designed pairs of primers for 8 human tissues. Actin cDNA was amplified as control.

After applying the new conditions, testis was the only tissue where products were constantly and clearly amplified. To clarify the transcription pattern of ATXN3L1, a better optimization of the amplification protocol must be performed.

## Chapter 4

# Discussion

---



## Chapter 4

# Discussion

---

Retrotransposition and retrogenes are gaining increasing attention as recent studies have shown that they may play an important role in genome evolution and in the formation of new genes.<sup>1,72</sup>

This project allowed us to investigate the evolution of *ataxin-3* paralogues formed by retrotransposition events, and to shed light on the functional relevance of these copies. Our results showed that *ATXN3* gene has been on the origin of many different retrocopies along the mammalian lineage. Two of these retrocopies, *ATXN3L1* and *ATXN3L2*, were found to have orthologues in several species from the primate lineage. The remaining copies, all diverging from each other, resulted from independent-origin retrotransposition events occurred after the speciation of several mammalian species. These observations are compatible with studies performed by Pan and Zhang (2007 and 2009) where they explain that a burst of young and independent retrogenes has recently arisen in several species branches, mainly in mammalian genomes.<sup>2,73</sup> The key factor underlying this important active process, that has been shaping the dynamics of mammalian genomes, might have been the recruitment of specific L1 retrotransposons in mammals.<sup>2,73</sup>

In the case of *ATXN3L1* and *ATXN3L2*, the identification of several orthologues in the primate lineage suggested an ancestral origin for each one of these copies, before the split of several primate clades, with their evolution until these days. Our results suggested that the two retrocopies have independent origins, with *ATXN3L1* as the most ancient copy (63 MYA), followed by the birth of *ATXN3L2*, about 35 MYA. Across their evolution, these copies had the chance to become functional or remain as non-relevant pseudogenes, having consequently more or less influence in the respective genomes of the considered species. Our results have shown some evidences that *ATXN3L1* is being conserved and its ORF kept in check. In addition, our results suggested that *ATXN3L1* is more likely to be transcribed than *ATXN3L2*. This difference may have resulted from the genomic background where copies were inserted. A study by Emerson et al. (2004) claims that retrogenes are not randomly located on chromosomes, but that genes are more likely to be retroposed into and/or out-of the X chromosome in mammals.<sup>74</sup> The authors demonstrated that, during evolution of human and

mouse genomes, the mammalian X chromosome has generated and recruited a disproportionately high number of functional retroposed genes, whereas the autosomes experienced lower gene turnover.<sup>8, 74</sup> Therefore, it is interesting to notice that the most conserved *ATXN3* retrocopy, *ATXN3L1*, is located on the X chromosome, whereas *ATXN3L2* (suggested by our results to be, most likely, a processed pseudogene) resides in an autosome.

As it was previously described, to become fixed in a given species genome, retrocopies must arise in the germ line and the gene source must be expressed at that time.<sup>1</sup> *ATXN3* is ubiquitously present in human tissues and cell types however, the retropositions of *ATXN3L1* and *ATXN3L2* occurred, not in the human species, but in an ancestor of Haplorrhines and Catarrhines, respectively.<sup>25, 38</sup> Thus, to be feasible the occurrence of retrotransposition events from *ATXN3*, this gene must have been transcribed in the germ line of these ancestral species.

The analysis of the (CAG)<sub>n</sub> tract of *ATXN3*, among different primate species, is a key factor to understand the process of repeat instability over the primate lineage until reaching the expanded/pathogenic range in humans. In the parental gene, the repeat tract configuration has been maintained stable along the primate lineage, varying mainly on the repeat size, which increased over time. Most of the studies have described a range between 14-40 in humans, 14-20 in chimps, 8-11 in gorillas and 24-25 in orangutans.<sup>75</sup> According to Andrés et al. (2004), humans show, not only a higher mean number of repeats than other species, but also a higher variance and coefficient of variation, with a decreasing trend observed as more anciently diverged primates are analysed. Thus, a relationship has been suggested between variability levels and expansion potential. In addition, the authors hypothesized the existence of some balancing selection favouring the existence of different alleles in human populations that would maintain the high (CAG)<sub>n</sub> diversity currently observed in our species.<sup>75</sup>

As for *ATXN3L1*, the CAG repeat showed a poorly polymorphic and highly interrupted tract across the primate lineage, whereas *ATXN3L2* presented a pure (CAG)<sub>n</sub> in some species, and a polymorphic hexanucleotide repeat in humans and chimpanzees. This recent acquisition of a repetitive CGGCAG resulted, probably, from a CAG to CGG mutation followed by instability that encompassed the six bases instead of the CAG alone.

In addition, studies on the CAG region of both parental and *ATXN3* paralogues can be important to gain insight into the involvement of these repetitive tracts in both pathological and non-pathological states of the protein. In fact, Schaefer et al. (2012) have claimed that (CAG)<sub>n</sub> tracts are not randomly distributed in the genome.<sup>76</sup> Also, recent studies have

demonstrated that the polyQ tract may actually have a function in repeat-containing proteins. In fact, we can observe a biased involvement of these proteins towards functions related to transcriptional regulation and nuclear localization.<sup>77, 78</sup> Moreover, recent studies have suggested the involvement of polyQ repeats in protein-protein interaction networks, by stabilizing structural changes to facilitate the interaction between own coiled-coil regions with coiled-coil regions of other proteins (Figure 14).<sup>76</sup> These facts suggest that (CAG)<sub>n</sub> tracts can actually be subjected to selective pressures, in order to maintain a possible functional role in the protein context.

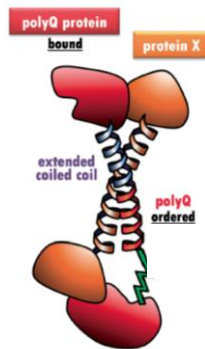


Figure 14 – Possible role for the polyQ stretch in protein-protein interactions: upon interaction with a protein partner, the polyQ region adopts a coiled-coil structure. (adapted from Schaefer, 2012)<sup>76</sup>

Analysing ORF predictions for both *ATXN3L1* and *ATXN3L2*, it was evident that the first copy, although more ancient, maintained its intact putative coding sequence. In opposition, the most recent retrocopy *ATXN3L2* displayed a disrupted ORF, with maximum length fragments of 200 bp predicted as alternative ORFs. These facts come to reinforce that both copies have accumulated genetic differences relatively to the parental gene along their evolution; nevertheless, *ATXN3L1* seems to be under selective constraints that might be keeping its sequence in check whereas *ATXN3L2* has accumulated premature stop codons that interrupted its original sequence. Yet, the short ORF of *ATXN3L2* may be transcribed into mRNA and play a trans-regulatory role in the expression of the parental gene. Indeed, it was recently found that some mammalian retropseudogenes have evolved the capacity to encode small interference RNAs, important for the regulation of their parental source genes.

Concerning the analysis of *ATXN3L1* expression in human tissues, our results have demonstrated that this paralogue is expressed, at least, in human testis. Results for other human tissues, such as brain, prostate, colon, liver and ovary, were less clear and, thus, need to be confirmed with further analyses. By comparing our results with expression data available in NCBI (EST profile) and in the Human Protein Atlas (<http://www.proteinatlas.org/>, updated at

11/11/2011), data is not consensual. Figure 15 shows a summary of the *ATXN3L1* expression levels, and Figure A7 of the appendix section a more complete version of these data.

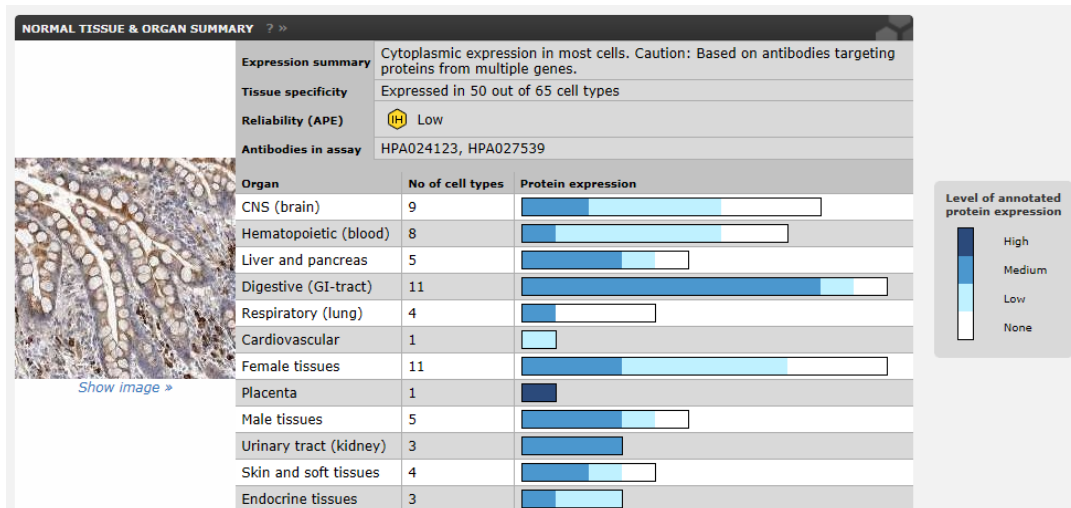


Figure 15 – *ATXN3L1* expression levels in different human organs/tissues. (data from the Human Protein Atlas: <http://www.proteinatlas.org/ENSG00000123594>)

It is important, however, to notice that expression patterns annotated in The Human Protein Atlas are based on antibody assays classified as “Low Reliability” in the case of *ATXN3L1*. On the other hand, *ATXN3L1* expression patterns consulted in the NCBI – UniGene – EST Profile are illustrated in Figure A8 of the appendix section. Here, the expression levels were measured by EST counts in unbiased cDNA libraries. Testis and brain are the only human body tissues where *ATXN3L1* is shown to be expressed, with higher levels of expression in testis. These data are, in fact, more similar to our results, but no further conclusions can be taken before confirming the obtained results

By comparing the principal domains of *ATXN3* and *ATXN3L1* proteins, we observed a conserved Josephin domain (catalytic domain), NES77, NES141 and UIM1, whereas less similarity was found for the remaining regions NLS273, UIM2 and UIM3. As the catalytic domain is the most conserved one (with intact catalytic sites: C14, H119 and N134), it can be hypothesized that *ATXN3L1* is likely to exert an ubiquitin protease activity as the parental *ATXN3*. The two other important residues in UIMs (serine-232 and serine-260 in UIM1 and



UIM2), as well as the UIM1 itself, are also conserved in ATXN3L1 suggesting that ubiquitin interactions can also be carried out by this retrocopy. However, possible functions of ATXN3L1 cannot be speculated based only on its conserved domains. More functional studies must be performed, for example, in order to explore common binding partners to the parental protein.



## Chapter 5

# Final remarks and future perspectives

---



## Chapter 5

# Final remarks and future perspectives

---

The general aims traced for this project were accomplished, but still, there are techniques here developed that can be optimized, as well as the enlargement of sample sizes in order to enrich and reinforce the obtained results.

From our study, it was possible to conclude that the analysis of paralogue genes involved in neurodegenerative diseases can be of great importance to gain insight into the mechanisms of pathogenesis. Namely, on Machado-Joseph disease, the ATXN3L1 was observed to be transcribed in human tissues; expression patterns will be assessed next through the optimization of specific amplifications of ATXN3L1 cDNA. At the protein level, ATXN3L1 possesses a highly conserved catalytic domain. Next, to detect the endogenous protein in human tissues, an anti-ATXN3L1 specific antibody can be generated to perform Western blot and immunohistochemistry. Additionally, future studies on the protein characterization will allow us to evaluate if this retrocopy has evolved towards functional diversification (neofunctionalization or neofunctionalization) relatively to the parental protein. To test if ATXN3L1 shares some of the most important molecular interactors (common binding partners) with ATXN3, one could detect common binding partners by co-immunoprecipitation with the transcriptional coactivators and DNA repair proteins that interact with the parental ataxin-3.

Further studies can be made, also, to confirm if the small ATXN3L2 ORFs predicted in this study are transcribed or not. This can be done by following the same approach we used here for ATXN3L1: the design of specific primers to amplify concatenated regions corresponding to ATXN3 exons, in gene-specific nucleotides with highly conserved sites across species within orthologues. If the transcription of ATXN3L2 occurs, one could explore the possibility of them being encoding small interference RNAs and acting as trans-regulatory factors of ATXN3.



## References

---





## References

---

1. Kaessmann H, Vinckenbosch N, Long M. RNA-based gene duplication: mechanistic and evolutionary insights. *Nat Rev Genet.* Jan 2009;10(1):19-31.
2. Pan D, Zhang L. Burst of young retrogenes and independent retrogene formation in mammals. *PLoS One.* 2009;4(3):e5040.
3. Ostertag EM, Kazazian HH, Jr. Biology of mammalian L1 retrotransposons. *Annu Rev Genet.* 2001;35:501-538.
4. Kazazian HH, Jr., Goodier JL. LINE drive, retrotransposition and genome instability. *Cell.* Aug 9 2002;110(3):277-280.
5. Vinckenbosch N, Dupanloup I, Kaessmann H. Evolutionary fate of retroposed gene copies in the human genome. *Proc Natl Acad Sci U S A.* Feb 28 2006;103(9):3220-3225.
6. Cordaux R, Batzer MA. The impact of retrotransposons on human genome evolution. *Nat Rev Genet.* Oct 2009;10(10):691-703.
7. Esnault C, Maestre J, Heidmann T. Human LINE retrotransposons generate processed pseudogenes. *Nat Genet.* Apr 2000;24(4):363-367.
8. Betran E, Thornton K, Long M. Retroposed new genes out of the X in *Drosophila*. *Genome Res.* Dec 2002;12(12):1854-1859.
9. Burki F, Kaessmann H. Birth and adaptive evolution of a hominoid gene that supports high neurotransmitter flux. *Nat Genet.* Oct 2004;36(10):1061-1063.
10. Carlson KM, Andresen JM, Orr HT. Emerging pathogenic pathways in the spinocerebellar ataxias. *Curr Opin Genet Dev.* Jun 2009;19(3):247-253.
11. Teive HA. Spinocerebellar ataxias. *Arq Neuropsiquiatr.* Dec 2009;67(4):1133-1142.
12. Soong BW, Paulson HL. Spinocerebellar ataxias: an update. *Curr Opin Neurol.* Aug 2007;20(4):438-446.
13. Bettencourt C, Lima M. Machado-Joseph Disease: from first descriptions to new perspectives. *Orphanet J Rare Dis.* 2011;6:35.
14. Kieling C, Morales Saute JA, Jardim LB. When ataxia is not just ataxia. *Nat Clin Pract Neurol.* May 2007;3(5):E2.
15. Coutinho P, Andrade C. Autosomal dominant system degeneration in Portuguese families of the Azores Islands. A new genetic disorder involving cerebellar, pyramidal, extrapyramidal and spinal cord motor functions. *Neurology.* Jul 1978;28(7):703-709.
16. Jardim LB, Silveira I, Pereira ML, et al. A survey of spinocerebellar ataxia in South Brazil - 66 new cases with Machado-Joseph disease, SCA7, SCA8, or unidentified disease-causing mutations. *J Neurol.* Oct 2001;248(10):870-876.
17. Vale J, Bugalho P, Silveira I, Sequeiros J, Guimaraes J, Coutinho P. Autosomal dominant cerebellar ataxia: frequency analysis and clinical characterization of 45 families from Portugal. *Eur J Neurol.* Jan 2010;17(1):124-128.
18. Zhao Y, Tan EK, Law HY, Yoon CS, Wong MC, Ng I. Prevalence and ethnic differences of autosomal-dominant cerebellar ataxia in Singapore. *Clin Genet.* Dec 2002;62(6):478-481.
19. Tang B, Liu C, Shen L, et al. Frequency of SCA1, SCA2, SCA3/MJD, SCA6, SCA7, and DRPLA CAG trinucleotide repeat expansion in patients with hereditary spinocerebellar ataxia from Chinese kindreds. *Arch Neurol.* Apr 2000;57(4):540-544.
20. van de Warrenburg BP, Sinke RJ, Verschuuren-Bemelmans CC, et al. Spinocerebellar ataxias in the Netherlands: prevalence and age at onset variance analysis. *Neurology.* Mar 12 2002;58(5):702-708.
21. Schols L, Amoiridis G, Buttner T, Przuntek H, Epplen JT, Riess O. Autosomal dominant cerebellar ataxia: phenotypic differences in genetically defined subtypes? *Ann Neurol.* Dec 1997;42(6):924-932.
22. Maruyama H, Izumi Y, Morino H, et al. Difference in disease-free survival curve and regional distribution according to subtype of spinocerebellar ataxia: a study of 1,286 Japanese patients. *Am J Med Genet.* Jul 8 2002;114(5):578-583.
23. Martins S, Calafell F, Gaspar C, et al. Asian origin for the worldwide-spread mutational event in Machado-Joseph disease. *Arch Neurol.* Oct 2007;64(10):1502-1508.
24. Takiyama Y, Nishizawa M, Tanaka H, et al. The gene for Machado-Joseph disease maps to human chromosome 14q. *Nat Genet.* Jul 1993;4(3):300-304.

25. Ichikawa Y, Goto J, Hattori M, et al. The genomic structure and expression of MJD, the Machado-Joseph disease gene. *J Hum Genet.* 2001;46(7):413-422.
26. Bettencourt C, Santos C, Montiel R, et al. Increased transcript diversity: novel splicing variants of Machado-Joseph disease gene (ATXN3). *Neurogenetics.* May 2010;11(2):193-202.
27. Kawaguchi Y, Okamoto T, Taniwaki M, et al. CAG expansions in a novel gene for Machado-Joseph disease at chromosome 14q32.1. *Nat Genet.* Nov 1994;8(3):221-228.
28. Wang JL, Xiao B, Cui XX, et al. Analysis of SCA2 and SCA3/MJD repeats in Parkinson's disease in mainland China: genetic, clinical, and positron emission tomography findings. *Mov Disord.* Oct 15 2009;24(13):2007-2011.
29. Gilman S. The spinocerebellar ataxias. *Clin Neuropharmacol.* Nov-Dec 2000;23(6):296-303.
30. Maruyama H, Nakamura S, Matsuyama Z, et al. Molecular features of the CAG repeats and clinical manifestation of Machado-Joseph disease. *Hum Mol Genet.* May 1995;4(5):807-812.
31. van de Warrenburg BP, Hendriks H, Durr A, et al. Age at onset variance analysis in spinocerebellar ataxias: a study in a Dutch-French cohort. *Ann Neurol.* Apr 2005;57(4):505-512.
32. Kawakami H, Maruyama H, Nakamura S, et al. Unique features of the CAG repeats in Machado-Joseph disease. *Nat Genet.* Apr 1995;9(4):344-345.
33. Maciel P, Gaspar C, DeStefano AL, et al. Correlation between CAG repeat length and clinical features in Machado-Joseph disease. *Am J Hum Genet.* Jul 1995;57(1):54-61.
34. Carvalho DR, La Rocque-Ferreira A, Rizzo IM, Imamura EU, Speck-Martins CE. Homozygosity enhances severity in spinocerebellar ataxia type 3. *Pediatr Neurol.* Apr 2008;38(4):296-299.
35. Panigrahi GB, Lau R, Montgomery SE, Leonard MR, Pearson CE. Slipped (CTG)\*(CAG) repeats can be correctly repaired, escape repair or undergo error-prone repair. *Nat Struct Mol Biol.* Aug 2005;12(8):654-662.
36. Martins S, Calafell F, Wong VC, Sequeiros J, Amorim A. A multistep mutation mechanism drives the evolution of the CAG repeat at MJD/SCA3 locus. *Eur J Hum Genet.* Aug 2006;14(8):932-940.
37. Djian P, Hancock JM, Chana HS. Codon repeats in genes associated with human diseases: fewer repeats in the genes of nonhuman primates and nucleotide substitutions concentrated at the sites of reiteration. *Proc Natl Acad Sci U S A.* Jan 9 1996;93(1):417-421.
38. Trottier Y, Cancel G, An-Gourfinkel I, et al. Heterogeneous intracellular localization and expression of ataxin-3. *Neurobiol Dis.* Nov 1998;5(5):335-347.
39. Bilen J, Bonini NM. Genome-wide screen for modifiers of ataxin-3 neurodegeneration in *Drosophila*. *PLoS Genet.* Oct 2007;3(10):1950-1964.
40. Macedo-Ribeiro S, Cortes L, Maciel P, Carvalho AL. Nucleocytoplasmic shuttling activity of ataxin-3. *PLoS One.* 2009;4(6):e5834.
41. Paulson HL, Das SS, Crino PB, et al. Machado-Joseph disease gene product is a cytoplasmic protein widely expressed in brain. *Ann Neurol.* Apr 1997;41(4):453-462.
42. Bichelmeier U, Schmidt T, Hubener J, et al. Nuclear localization of ataxin-3 is required for the manifestation of symptoms in SCA3: in vivo evidence. *J Neurosci.* Jul 11 2007;27(28):7418-7428.
43. Chai Y, Koppenhafer SL, Shoesmith SJ, Perez MK, Paulson HL. Evidence for proteasome involvement in polyglutamine disease: localization to nuclear inclusions in SCA3/MJD and suppression of polyglutamine aggregation in vitro. *Hum Mol Genet.* Apr 1999;8(4):673-682.
44. Mukhopadhyay D, Riezman H. Proteasome-independent functions of ubiquitin in endocytosis and signaling. *Science.* Jan 12 2007;315(5809):201-205.
45. Todi SV, Laco MN, Winborn BJ, Travis SM, Wen HM, Paulson HL. Cellular turnover of the polyglutamine disease protein ataxin-3 is regulated by its catalytic activity. *J Biol Chem.* Oct 5 2007;282(40):29348-29358.
46. Chou AH, Yeh TH, Ouyang P, Chen YL, Chen SY, Wang HL. Polyglutamine-expanded ataxin-3 causes cerebellar dysfunction of SCA3 transgenic mice by inducing transcriptional dysregulation. *Neurobiol Dis.* Jul 2008;31(1):89-101.

47. Nicastro G, Todi SV, Karaca E, Bonvin AM, Paulson HL, Pastore A. Understanding the role of the Josephin domain in the PolyUb binding and cleavage properties of ataxin-3. *PLoS One*. 2010;5(8):e12430.
48. Masino L, Musi V, Menon RP, et al. Domain architecture of the polyglutamine protein ataxin-3: a globular domain followed by a flexible tail. *FEBS Lett*. Aug 14 2003;549(1-3):21-25.
49. Harris GM, Dodelzon K, Gong L, Gonzalez-Alegre P, Paulson HL. Splice isoforms of the polyglutamine disease protein ataxin-3 exhibit similar enzymatic yet different aggregation properties. *PLoS One*. 2010;5(10):e13695.
50. Scheel H, Tomiuk S, Hofmann K. Elucidation of ataxin-3 and ataxin-7 function by integrative bioinformatics. *Hum Mol Genet*. Nov 1 2003;12(21):2845-2852.
51. Mueller T, Breuer P, Schmitt I, Walter J, Evert BO, Wullner U. CK2-dependent phosphorylation determines cellular localization and stability of ataxin-3. *Hum Mol Genet*. Sep 1 2009;18(17):3334-3343.
52. Weeks SD, Grasty KC, Hernandez-Cuebas L, Loll PJ. Crystal structure of a Josephin-ubiquitin complex: evolutionary restraints on ataxin-3 deubiquitinating activity. *J Biol Chem*. Feb 11 2011;286(6):4555-4565.
53. Matos CA, de Macedo-Ribeiro S, Carvalho AL. Polyglutamine diseases: the special case of ataxin-3 and Machado-Joseph disease. *Prog Neurobiol*. Sep 15 2011;95(1):26-48.
54. Ferrigno P, Silver PA. Polyglutamine expansions: proteolysis, chaperones, and the dangers of promiscuity. *Neuron*. Apr 2000;26(1):9-12.
55. Muchowski PJ, Schaffar G, Sittler A, Wanker EE, Hayer-Hartl MK, Hartl FU. Hsp70 and hsp40 chaperones can inhibit self-assembly of polyglutamine proteins into amyloid-like fibrils. *Proc Natl Acad Sci U S A*. Jul 5 2000;97(14):7841-7846.
56. Chai Y, Wu L, Griffin JD, Paulson HL. The role of protein composition in specifying nuclear inclusion formation in polyglutamine disease. *J Biol Chem*. Nov 30 2001;276(48):44889-44897.
57. Li F, Macfarlan T, Pittman RN, Chakravarti D. Ataxin-3 is a histone-binding protein with two independent transcriptional corepressor activities. *J Biol Chem*. Nov 22 2002;277(47):45004-45012.
58. Evert BO, Araujo J, Vieira-Saecker AM, et al. Ataxin-3 represses transcription via chromatin binding, interaction with histone deacetylase 3, and histone deacetylation. *J Neurosci*. Nov 1 2006;26(44):11474-11486.
59. Gunawardena S, Her LS, Brusch RG, et al. Disruption of axonal transport by loss of huntingtin or expression of pathogenic polyQ proteins in *Drosophila*. *Neuron*. Sep 25 2003;40(1):25-40.
60. Donaldson KM, Li W, Ching KA, Batalov S, Tsai CC, Joazeiro CA. Ubiquitin-mediated sequestration of normal cellular proteins into polyglutamine aggregates. *Proc Natl Acad Sci U S A*. Jul 22 2003;100(15):8892-8897.
61. Paulson HL, Perez MK, Trotter Y, et al. Intranuclear inclusions of expanded polyglutamine protein in spinocerebellar ataxia type 3. *Neuron*. Aug 1997;19(2):333-344.
62. Goti D, Katzen SM, Mez J, et al. A mutant ataxin-3 putative-cleavage fragment in brains of Machado-Joseph disease patients and transgenic mice is cytotoxic above a critical concentration. *J Neurosci*. Nov 10 2004;24(45):10266-10279.
63. Yoshizawa T, Yamagishi Y, Koseki N, et al. Cell cycle arrest enhances the in vitro cellular toxicity of the truncated Machado-Joseph disease gene product with an expanded polyglutamine stretch. *Hum Mol Genet*. Jan 1 2000;9(1):69-78.
64. Haacke A, Broadley SA, Boteva R, Tzvetkov N, Hartl FU, Breuer P. Proteolytic cleavage of polyglutamine-expanded ataxin-3 is critical for aggregation and sequestration of non-expanded ataxin-3. *Hum Mol Genet*. Feb 15 2006;15(4):555-568.
65. Wellington CL, Ellerby LM, Hackam AS, et al. Caspase cleavage of gene products associated with triplet expansion disorders generates truncated fragments containing the polyglutamine tract. *J Biol Chem*. Apr 10 1998;273(15):9158-9167.
66. Hubener J, Vauti F, Funke C, et al. N-terminal ataxin-3 causes neurological symptoms with inclusions, endoplasmic reticulum stress and ribosomal dislocation. *Brain*. Jul 2011;134(Pt 7):1925-1942.
67. Costa Mdo C, Paulson HL. Toward understanding Machado-Joseph disease. *Prog Neurobiol*. May 2012;97(2):239-257.
68. Fujigasaki H, Uchihara T, Koyano S, et al. Ataxin-3 is translocated into the nucleus for the formation of intranuclear inclusions in normal and Machado-Joseph disease brains. *Exp Neurol*. Oct 2000;165(2):248-256.
69. Warrick JM, Morabito LM, Bilen J, et al. Ataxin-3 suppresses polyglutamine neurodegeneration in *Drosophila* by a ubiquitin-associated mechanism. *Mol Cell*. Apr 1 2005;18(1):37-48.

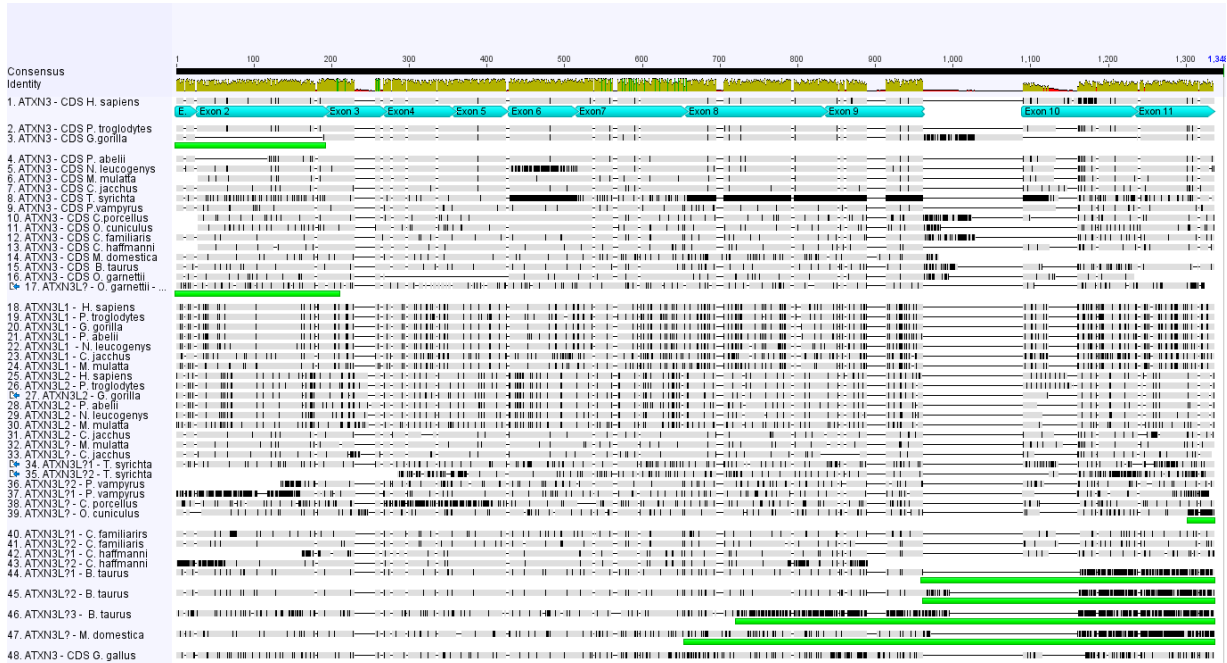
70. Crespo-Barreto J, Fryer JD, Shaw CA, Orr HT, Zoghbi HY. Partial loss of ataxin-1 function contributes to transcriptional dysregulation in spinocerebellar ataxia type 1 pathogenesis. *PLoS Genet.* Jul 2010;6(7):e1001021.
71. Bowman AB, Lam YC, Jafar-Nejad P, et al. Duplication of Atxn1l suppresses SCA1 neuropathology by decreasing incorporation of polyglutamine-expanded ataxin-1 into native complexes. *Nat Genet.* Mar 2007;39(3):373-379.
72. Kaessmann H. Genetics. More than just a copy. *Science.* Aug 21 2009;325(5943):958-959.
73. Pan D, Zhang L. Quantifying the major mechanisms of recent gene duplications in the human and mouse genomes: a novel strategy to estimate gene duplication rates. *Genome Biol.* 2007;8(8):R158.
74. Emerson JJ, Kaessmann H, Betran E, Long M. Extensive gene traffic on the mammalian X chromosome. *Science.* Jan 23 2004;303(5657):537-540.
75. Andres AM, Soldevila M, Lao O, et al. Comparative genetics of functional trinucleotide tandem repeats in humans and apes. *J Mol Evol.* Sep 2004;59(3):329-339.
76. Schaefer MH, Wanker EE, Andrade-Navarro MA. Evolution and function of CAG/polyglutamine repeats in protein-protein interaction networks. *Nucleic Acids Res.* May 1 2012;40(10):4273-4287.
77. Harrison PM. Exhaustive assignment of compositional bias reveals universally prevalent biased regions: analysis of functional associations in human and *Drosophila*. *BMC Bioinformatics.* 2006;7:441.
78. Alba MM, Guigo R. Comparative analysis of amino acid repeats in rodents and humans. *Genome Res.* Apr 2004;14(4):549-554.
79. Jobling MA, Hurles, ME, Tyler-Smith C. *Human Evolutionary Genetics – Origins, Peoples & Disease.* Garland Publishing (2004)

# Appendix

---



# Appendix



**Figure A1** – Alignment of all sequences searched on Ensembl, NCBI and UCSC Genome Browser databases. From top to bottom: *Gallus gallus* as outgroup; ATXN3 coding sequences for primates; ATXN3 coding sequences for non-primate mammals; ATXN3L1 for primates; ATXN3L2 for primates; ATXN3L3 of primates and non-primate mammals. The human exons present in transcript ATXN3-001 were annotated in the ATXN3 coding sequence to serve as reference. Green annotations correspond to sequence portions added to those sequences obtained from trace archives.







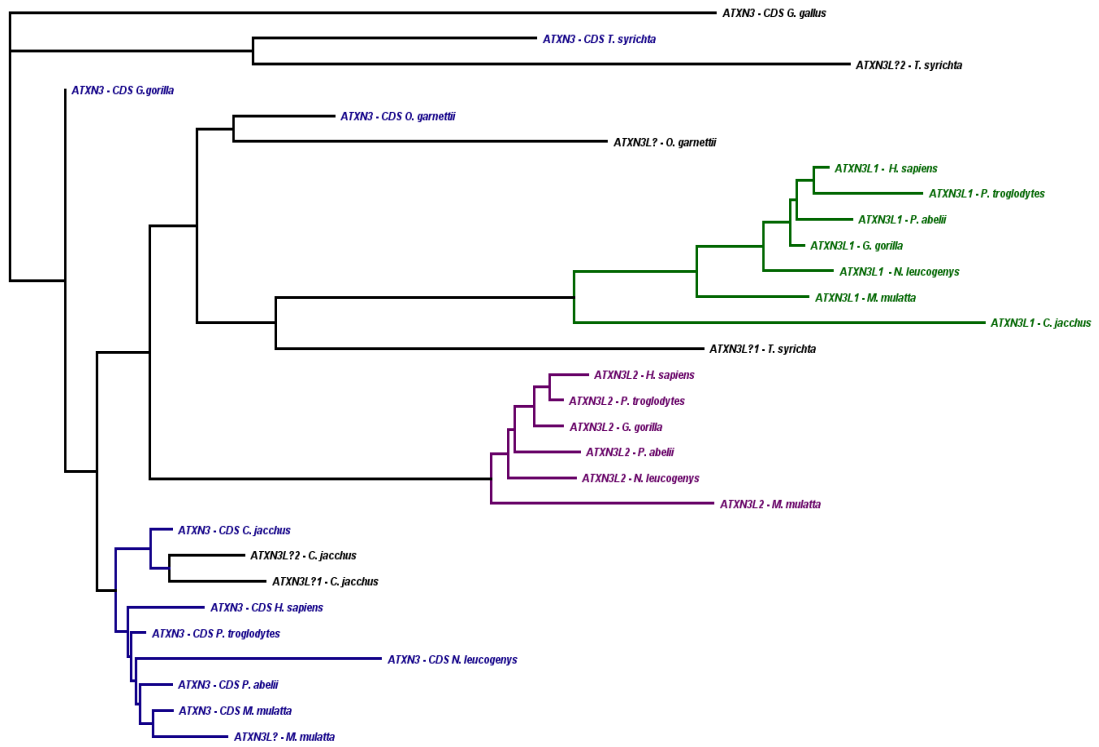


Figure A2 –Phylogenetic tree of *ATXN3*, *ATXN3L1* and *ATXN3L2* for primates, without the GAG tract, based on the genetic distances presented on Table A2. *Gallus gallus* was used as outgroup.

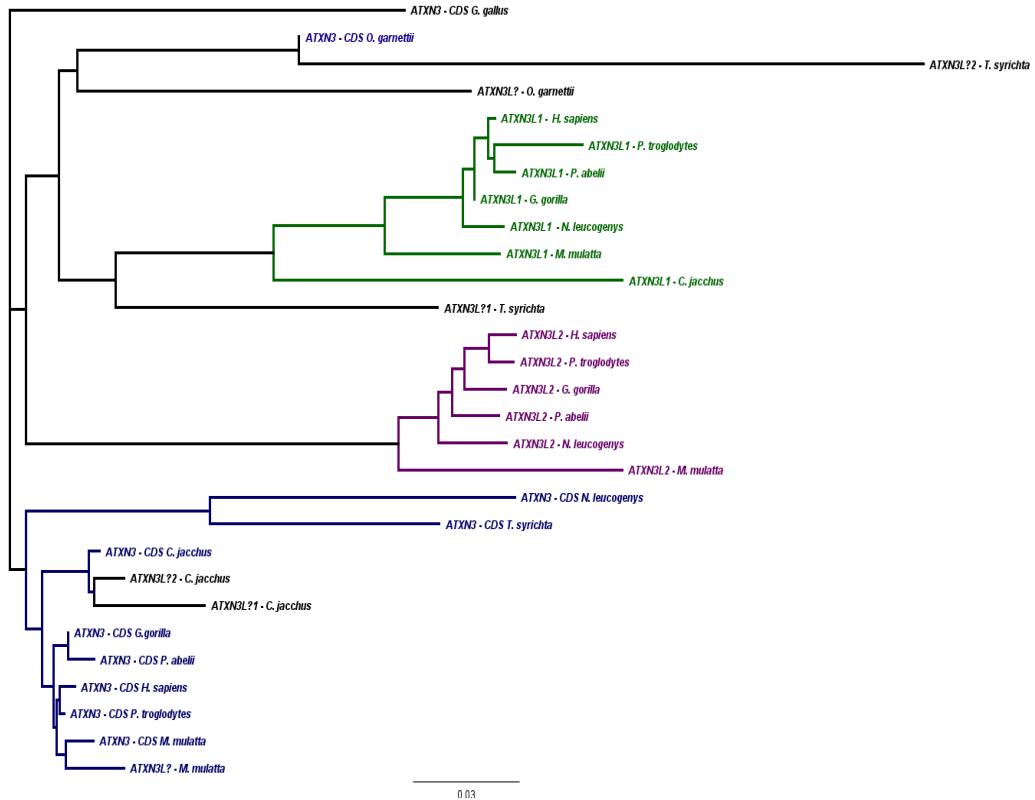


Figure A3 – Phylogenetic tree of *ATXN3*, *ATXN3L1* and *ATXN3L2* Josephin domain for primates, based on the genetic distances presented on Table A3. *Gallus gallus* was used as outgroup.

















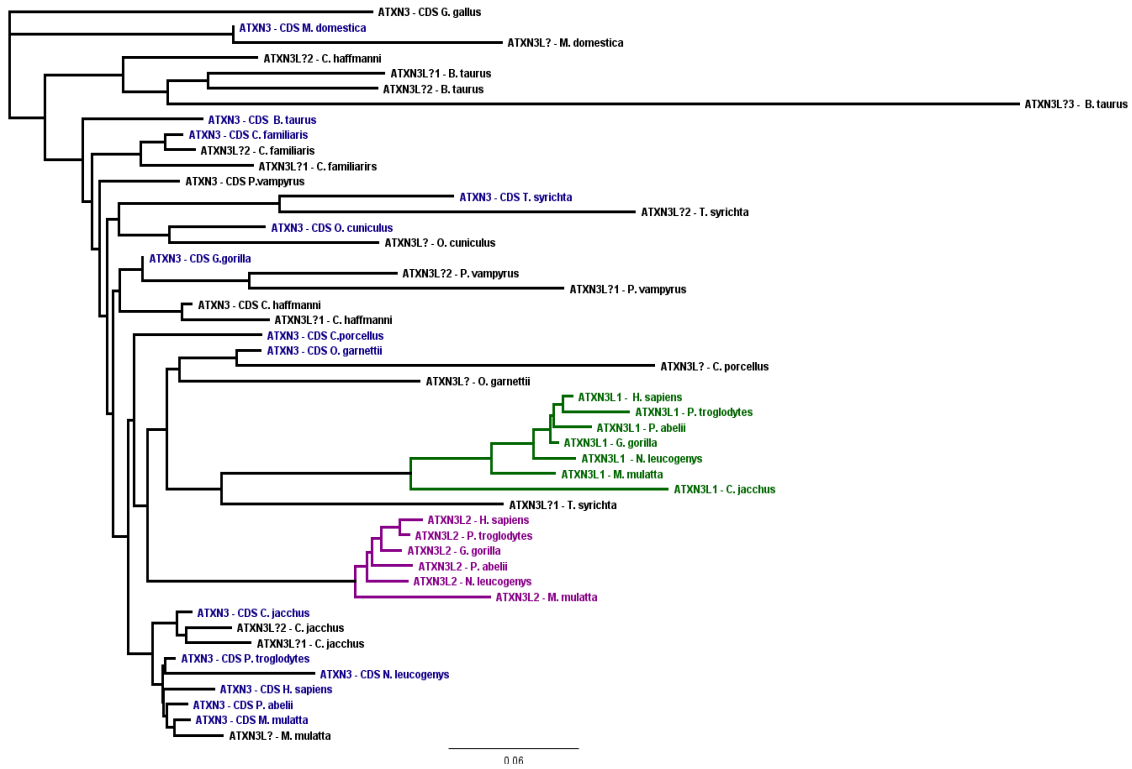


Figure A4 – Phylogenetic tree of *ATXN3*, *ATXN3L1* and *ATXN3L2* for mammals, without the GAG tract, based on the genetic distances presented on Table A5. *Gallus gallus* was used as outgroup.

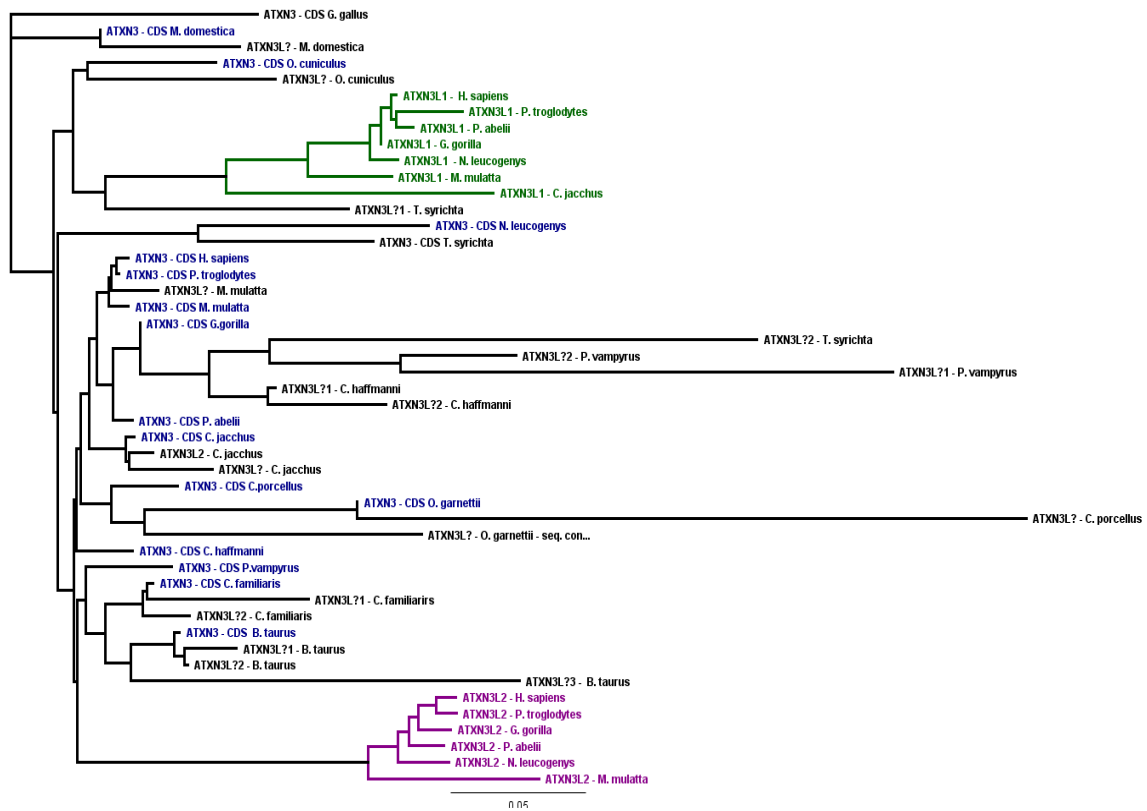


Figure A5 – Phylogenetic tree of *ATXN3*, *ATXN3L1* and *ATXN3L2* Josephin domain for mammals, based on the genetic distances presented on table A6. *Gallus gallus* was used as outgroup.

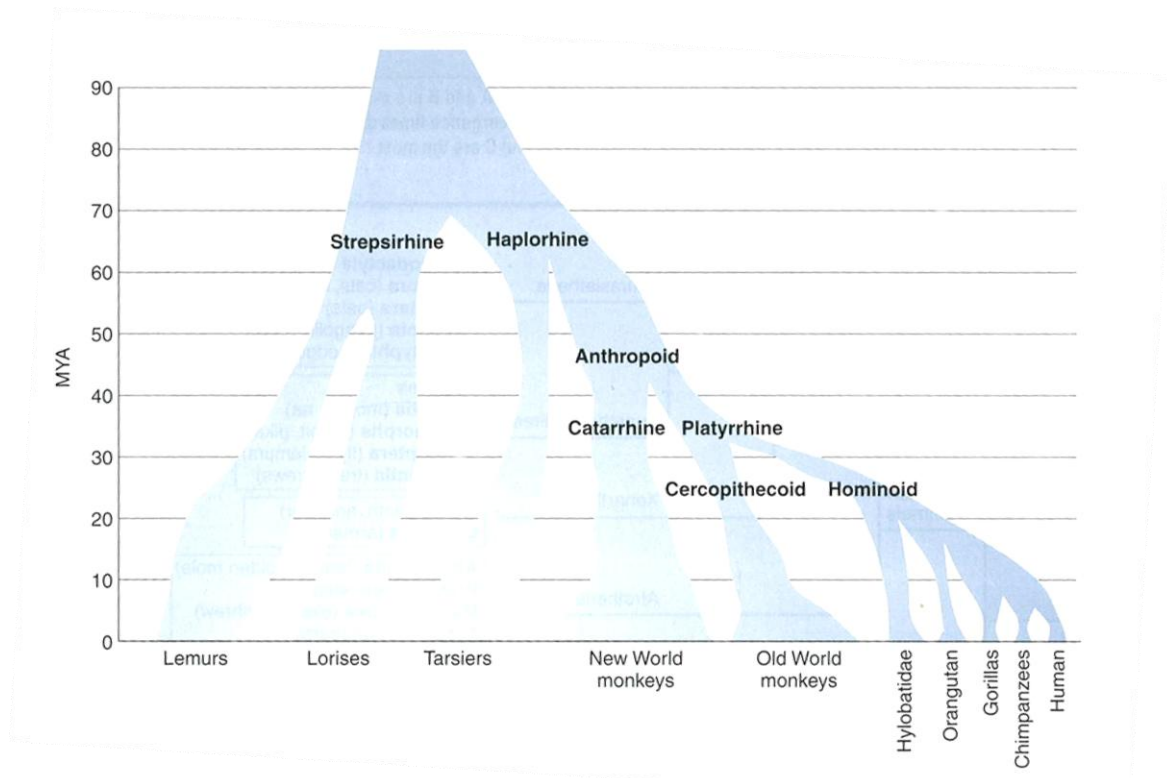


Figure A6 – Phylogeny of extant primate groups and corresponding date of branching splits.<sup>79</sup>

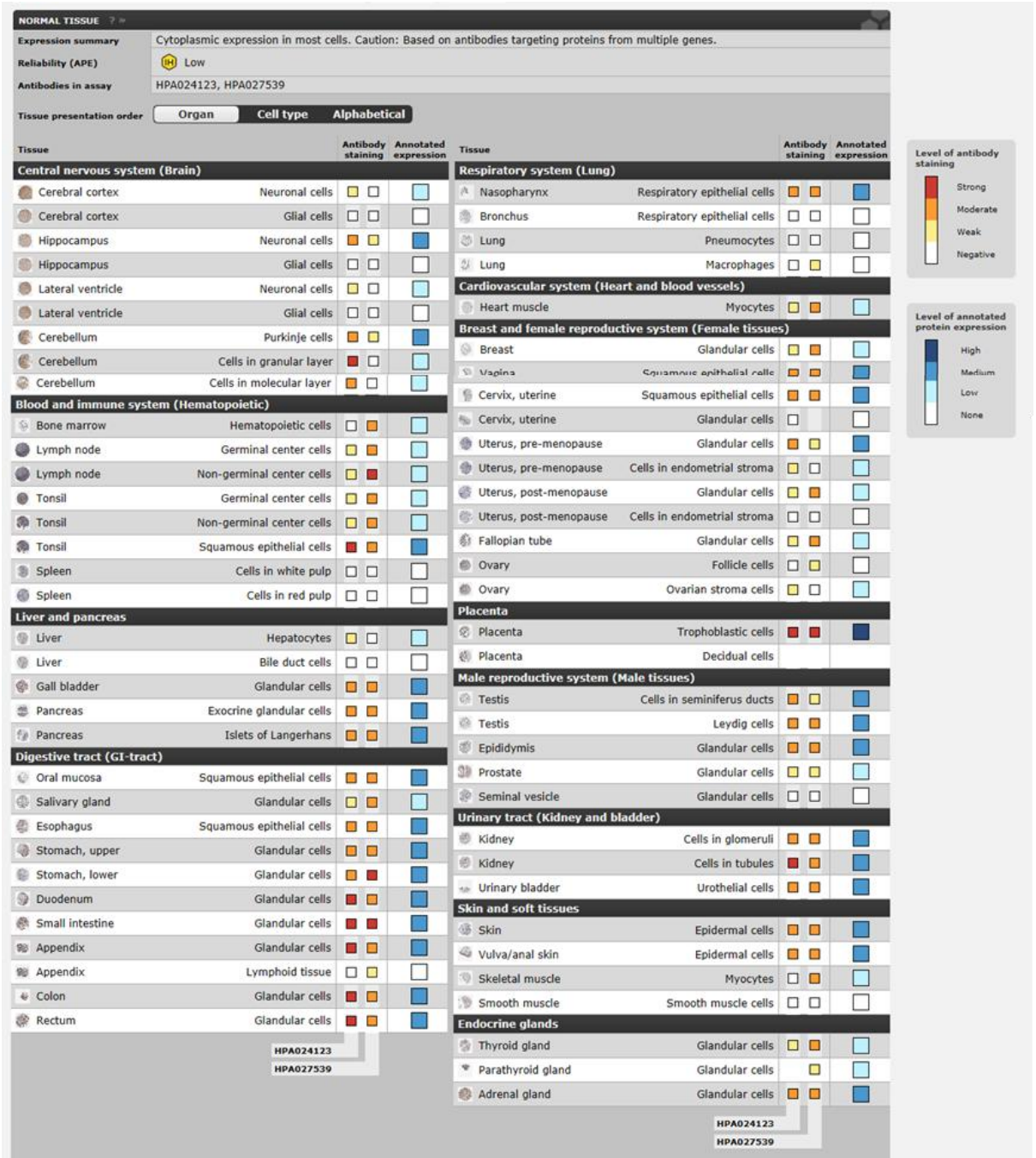
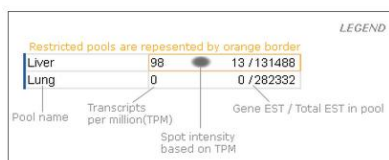


Figure A7 - ATXN3L1 expression levels in different human organs/tissues. (data from the Human Protein Atlas: <http://www.proteinatlas.org/ENSG00000123594>)

Hs.382641 - ATXN3L: Ataxin 3-like

**Breakdown by Body Sites**

	Hs.382641	
adipose tissue	0	0/12865
adrenal gland	0	0/32921
ascites	0	0/39829
bladder	0	0/29856
blood	0	0/122262
bone	0	0/71609
bone marrow	0	0/48711
brain	1	2/1092524
cervix	0	0/48469
connective tissue	0	0/149048
ear	0	0/16098
embryonic tissue	0	0/212847
esophagus	0	0/20152
eye	0	0/208810
heart	0	0/89512
intestine	0	0/232030
kidney	0	0/210738
larynx	0	0/23489
liver	0	0/205232
lung	0	0/334751
lymph	0	0/44292
lymph node	0	0/89697
mammary gland	0	0/151228
mouth	0	0/66139
muscle	0	0/106323
nerve	0	0/15526
ovary	0	0/101482
pancreas	0	0/213410
parathyroid	0	0/20579
pharynx	0	0/40767
prostate	0	0/189585
salivary gland	0	0/20264
skin	0	0/210718
spleen	0	0/53365
stomach	0	0/95775
testis	21	7/327305
thymus	0	0/79668
thyroid	0	0/46584
tonsil	0	0/17016
trachea	0	0/51769
umbilical cord	0	0/13764
uterus	0	0/232051
vascular	0	0/51637



**Figure A8 - ATXN3L1 expression levels in different human organs/tissues. (data from NCBI – UniGene – EST Profile data base: <http://www.ncbi.nlm.nih.gov/UniGene/ESTProfileViewer.cgi?uglis t=Hs.382641>)**

