



HAL
open science

Extensive recombination events and horizontal gene transfer shaped the *Legionella pneumophila* genomes

Laura Gomez-Valero, Christophe Rusniok, Sophie Jarraud, Benoit Vacherie, Zoé Rouy, Valérie Barbe, Claudine Médigue, Jerome Etienne, Carmen Buchrieser

► **To cite this version:**

Laura Gomez-Valero, Christophe Rusniok, Sophie Jarraud, Benoit Vacherie, Zoé Rouy, et al.. Extensive recombination events and horizontal gene transfer shaped the *Legionella pneumophila* genomes. BMC Genomics, 2011, 12 (1), pp.536. 10.1186/1471-2164-12-536 . pasteur-00642457

HAL Id: pasteur-00642457

<https://pasteur.hal.science/pasteur-00642457>

Submitted on 18 Nov 2011

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

RESEARCH ARTICLE

Open Access

Extensive recombination events and horizontal gene transfer shaped the *Legionella pneumophila* genomes

Laura Gomez-Valero^{1,2}, Christophe Rusniok^{1,2}, Sophie Jarraud^{3,4,5}, Benoit Vacherie⁶, Zoé Rouy^{6,7}, Valerie Barbe⁶, Claudine Medigue^{6,7}, Jerome Etienne^{3,4,5} and Carmen Buchrieser^{1,2*}

Abstract

Background: *Legionella pneumophila* is an intracellular pathogen of environmental protozoa. When humans inhale contaminated aerosols this bacterium may cause a severe pneumonia called Legionnaires' disease. Despite the abundance of dozens of *Legionella* species in aquatic reservoirs, the vast majority of human disease is caused by a single serogroup (Sg) of a single species, namely *L. pneumophila* Sg1. To get further insights into genome dynamics and evolution of Sg1 strains, we sequenced strains Lorraine and HL 0604 1035 (Sg1) and compared them to the available sequences of Sg1 strains Paris, Lens, Corby and Philadelphia, resulting in a comprehensive multigenome analysis.

Results: We show that *L. pneumophila* Sg1 has a highly conserved and syntenic core genome that comprises the many eukaryotic like proteins and a conserved repertoire of over 200 Dot/Icm type IV secreted substrates. However, recombination events and horizontal gene transfer are frequent. In particular the analyses of the distribution of nucleotide polymorphisms suggests that large chromosomal fragments of over 200 kbs are exchanged between *L. pneumophila* strains and contribute to the genome dynamics in the natural population. The many secretion systems present might be implicated in exchange of these fragments by conjugal transfer. Plasmids also play a role in genome diversification and are exchanged among strains and circulate between different *Legionella* species.

Conclusion: Horizontal gene transfer among bacteria and from eukaryotes to *L. pneumophila* as well as recombination between strains allows different clones to evolve into predominant disease clones and others to replace them subsequently within relatively short periods of time.

Background

Legionella pneumophila is the etiologic agent of Legionnaires' disease, an atypical pneumonia, which is often fatal if not treated promptly. However, it is principally an environmental bacterium that inhabits fresh water reservoirs worldwide where it parasitizes within free-living protozoa but also survives in biofilms [1-3]. Since *L. pneumophila* does not spread from person-to-person, humans have been inconsequential for the evolution of this pathogen. Instead, the virulence strategies of *L. pneumophila* have been shaped by selective pressures in

aquatic ecosystems. Indeed, the co-evolution of *L. pneumophila* with fresh-water amoebae is reflected in its genome sequence. The analysis of two *L. pneumophila* genomes identified the presence of an unexpected high number and variety of eukaryotic-like proteins and proteins containing motifs mainly found in eukaryotes [4]. These proteins were predicted to interfere in different steps of the infectious cycle by mimicking functions of eukaryotic proteins [4]. For several of these eukaryotic like proteins it has been shown recently that they are secreted effectors that help *L. pneumophila* to subvert host functions to allow intracellular replication [5,6]. The possibility that *L. pneumophila* has acquired at least some of these genes through horizontal gene transfer from eukaryotes has been suggested by two studies [7,8].

* Correspondence: carmen.buchrieser@pasteur.fr

¹Institut Pasteur, Biologie des Bactéries Intracellulaires, 75724, Paris, France
Full list of author information is available at the end of the article

Plasticity is another specific feature of the *L. pneumophila* genomes as integrative plasmids, putative conjugation elements and genomic islands were identified. In addition to DNA interchange between different bacterial genera and even domains of life, horizontal gene transfer within the genus *Legionella* and within the species *L. pneumophila* has been reported. For example a 65-kb pathogenicity island described first in *L. pneumophila* strain Philadelphia [9] is present in several *L. pneumophila* strains and also in other *Legionella* species like *L. anisa* [10]. Another example is the particular lipopolysaccharide cluster of serogroup 1 strains that has been detected in *L. pneumophila* strains of different lineages and genetic backgrounds [10]. *L. pneumophila* has all necessary features for incorporating foreign DNA, as these bacteria are naturally competent and possess an intact recombination machinery [11,12]. These findings suggest that the *L. pneumophila* genomes are very dynamic and one would expect that horizontal gene transfer and recombination events play an important role in their evolution.

However, different analyses like early studies applying multilocus enzyme electrophoreses (MEE) supported a clonal population structure of *L. pneumophila* [13]. Two recent reports using genetic profiling based on six or three genetic loci, respectively concluded also that *L. pneumophila* shows a clonal populations structure [14,15] although the presence of few recombination events was not ruled out. Later the analysis of the *dotA*, *mip* and *rpoB* genes in different isolates suggested for the first time that recombination may play some role in *L. pneumophila* evolution [16-18] and a more in depth analysis using over 20 loci suggested that recombination events might be more frequent than was previously thought [19]. However, comparisons of these studies are difficult due to different sampling and different analysis methods used. Furthermore there may be a bias associated with some of the genes selected in these studies like intergenic spacer regions or genes under positive selection that may lead to artefactual effects in detecting recombination. To solve these problem efforts have been undertaken recently to homogenize the results obtained for different species to allow comparisons [20]. These authors report for *L. pneumophila* a low recombination rate like for the obligate pathogens *Bordetella pertussis* or *Bartonella henselae*. In contrast Coscolla and colleagues suggest a more important role for recombination at the intergenic level [21].

These different results and the fact that a globally distributed *L. pneumophila* clone implicated in Legionnaires' disease has been described [10] may suggest that the role of recombination is not relevant. However, the description of clonal complexes is not incompatible with high recombination rates. Transient clones may appear

within a recombining population [22], in particular if clones with high disease prevalence appear, as this seems to be the case for some *L. pneumophila* strains. These clones are often vastly over-sampled due to their clinical importance and show strong clonality. Thus, this may be correct for this subgroup, but it may not be representative for the population. Indeed when analyzing over 200 clinical and environmental *L. pneumophila* strains, significantly less diversity was found among the clinical isolates [23].

In this study we investigated the genome dynamics and evolution of the species *L. pneumophila* by analyzing horizontal gene transfer, mobile genetic elements and recombination on a genome-wide level. We undertook this analysis based on six complete genome sequences four of which are the previously published reference genomes of *L. pneumophila* Paris, Lens [4], Corby [24] and Philadelphia [25] and two that were sequenced in this study. The newly sequenced strains were selected according to epidemiological features that might be reflected in their genomes and should thus allow to study genome dynamics with respect to virulence. Strain Lorraine is rarely isolated from the environment but its prevalence in human disease is increasing considerably in the last years [26]. In contrast, *L. pneumophila* strain HL 0604 1035 has been frequently isolated from a hospital water system since over 10 years but has never caused disease. Analysis of these six strains identified a highly conserved and syntenic core genome and a diverse accessory genome. Furthermore, it showed that recombination events and horizontal gene transfer are frequent in *L. pneumophila*. Horizontal gene transfer from eukaryotes as well as recombination between strains were identified suggesting that *L. pneumophila* genomes are highly dynamic, a feature allowing different clones to evolve into predominant disease clones and others to replace them subsequently within relatively short periods of time.

Results and discussion

The *L. pneumophila* core genome comprises over 2400 conserved genes that are highly syntenic

To get comprehensive insight into the genetic basis, evolution and genome dynamics of *L. pneumophila* Sg1, the strains responsible for over 90% of disease worldwide, we analyzed six completely sequenced genomes. The strains selected are all of Sg1, have endemic and/or epidemic character (e.g. Paris, Lorraine or Philadelphia) were isolated in different countries (France, England, Spain, US) and in different years. Two strains were newly sequenced for this study (Lorraine and HL 0604 1035), the other four *L. pneumophila* genomes (Paris, Lens, Philadelphia, Corby) have been published previously [4,24,25]. The genomes of *L. pneumophila*

Table 1 General features of the 6 *L. pneumophila* strains analyzed

<i>L. pneumophila</i> strains	Philadelphia	Paris	Lens	Corby	HL06041035	Lorraine
Chromosome size (bp)	3397754	3503610	3345687	3576469	3492535	3467254
G+C content (%)	38.27	38.37	38.42	38.48	38.35	38.36
N° of genes	3031	3123	2980	3237	3132	3117
N° of protein coding genes	2999	3078	2921	3193	3079	3080
Pseudogenes	55	71	84	59	73	48
tRNA	43	43	43	44	43	44
16S/23S/5S	3/3/3	3/3/3	3/3/3	3/3/3	3/3/3	3/3/3
Average length CDS (nts)	1082.47	1000.85	1008.76	984.35	995.47	988.54
Average length ig (nts)	147.72	154	152.36	149.24	155.12	155.28
Coding density (%)	88.22	86.93	87.07	87.25	86.94	87.26
Plasmids	0	1	1	0	0	1

bp, base pairs; nts, nucleotides; CDS, coding sequence; ig, intergenic region

Lorraine and HL 0604 1035 consist each of a single circular chromosome of 3.4 Mb. Strain Lorraine also contains a plasmid. As shown in Table 1, the main features of the six *L. pneumophila* genomes analyzed (e.g. genome size, GC content and coding density), are highly conserved. The core genome of the six *L. pneumophila* genomes comprises 2434 genes, which represents about 80% of the predicted genes in each genome. Furthermore, the gene order is highly conserved as the 260 kb inversion in strain Lens with respect to the other strains is the only exception. When comparing the strains two by two, in average 90% of the genes are present in both strains (Figure 1). However, when determining the non-orthologous genes specific of each genome and not present in the remaining 5 strains, each strain contains between 136 (strain HL 0604 1035) and 222 (strain Corby) strain specific genes mainly encoded on mobile genetic elements. Taken together, the *L. pneumophila* genomes have a highly conserved and syntenic backbone and a highly dynamic accessory genome of about 300 genes each mainly formed by mobile genetic elements, genomic islands and genes of unknown function. The complete annotation of these six genomes is available in a new data base resource that we have set up, LegionellaScope https://www.genoscope.cns.fr/agc/microscope/about/collabprojects.php?P_id=27 and at the Institut Pasteur, LegioList <http://genolist.pasteur.fr/LegioList/>.

The species *L. pneumophila* has a highly conserved core genome

a) Most eukaryotic like proteins are conserved in all *L. pneumophila* genomes

The presence of proteins with high similarity to eukaryotic proteins or proteins with domains preferentially or only present in eukaryotic genomes are a particular feature of *L. pneumophila* [4]. However, the criteria for identifying these proteins were never clearly defined. To analyze their evolution and possible origin in depth we

have thus developed an automatic and systematic method to identify eukaryotic like proteins according to defined criteria. Previously we had identified eukaryotic like proteins in *L. pneumophila* as proteins with the

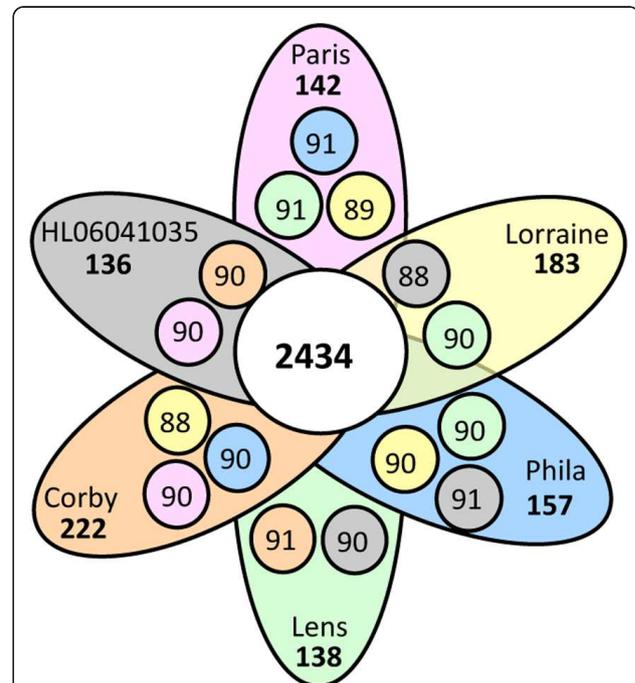


Figure 1 Shared and specific gene content of 6 *L. pneumophila* genomes. Each petal represents a genome with an associated color. The number in the center of the diagram represents the orthologous genes shared by all the genomes. The number inside of each individual petal corresponds to the specific genes of each genome with non-orthologous genes in any of the other genomes. The small circles inside of each petal represent the percentage of shared genes (total number divided by the number of genes in the smallest genome) between the genome of this petal and the genome represented by the color of the small circle. Yellow circle inside orange petal means that there are 88% of genes shared among Corby and Lorraine.

highest similarity score to eukaryotic proteins according to BLAST results, or by identifying eukaryotic domains [4]. However, due to constantly growing databases BLAST results are changing. Furthermore, recent analyses of amoeba-associated bacteria, in particular symbionts of amoeba have shown that they also contain eukaryotic like proteins, suggesting multiple origins of these proteins in prokaryotes [27]. To get a more complete picture of eukaryotic like proteins of *L. pneumophila* and also to include those proteins that might have been transferred independently to different amoeba associated bacteria we defined a eukaryotic like protein as i) a protein having a better normalized blast score against eukaryotic sequences than against prokaryotic ones and ii) a protein that did not show BLAST results against neither *Legionella* spp. nor other bacterial species for which resistance to amoeba infection has been demonstrated (see material and methods). Applying these criteria we identified 46 proteins with putative eukaryotic origin, of which 17 are described here for the first time (Table 2). Given the fact that these proteins were probably acquired by HGT one would expect high diversity in the repertoire. However, our analyses revealed a considerable conservation as more than 50% (26) are conserved in all six *L. pneumophila* strains, indicating an ancient transfer. Furthermore, they show 89-99% nucleotide identity, probably due to high selection pressure for their maintenance. Thus most of these proteins belong to the core genome, indicating that their acquisition has taken place before the speciation of *L. pneumophila*. These 26 proteins might have allowed a common *Legionella* ancestor to colonize an intracellular niche or to adapt better to the intercellular environment of a specific protozoan species leading to the evolution of the species *L. pneumophila*. Interestingly, 19 of these 26 proteins are also conserved in *L. longbeachae*, which might thus be those indispensable for intracellular replication of *Legionellae* (Table 2) [28].

b) Eukaryotic protein motifs are highly conserved among the *L. pneumophila* genomes

A second class of eukaryotic proteins of *L. pneumophila* is carrying domains predominantly present in eukaryotic proteins. To systematically identify these proteins we used the Interpro database comprising 10 different domain search programs [29]. This allowed to identify the *L. pneumophila* proteins carrying eukaryotic domains in the newly sequenced strains Lorraine and HL 0604 1035 as well as to identify previously not reported motifs. Similarly to the above described eukaryotic like proteins over half of the eukaryotic domain coding proteins are conserved in all six genomes and over 80% are conserved when two genomes are compared (e.g. 33 of the 39 proteins containing an eukaryotic motif in strain Lens are present also in strain Paris).

Moreover half of them share very high nucleotide identity of in average 98%-100% (Table 3) again suggesting high selection pressure to maintain them.

Our approach identified also new eukaryotic domains like spectrin repeats. The spectrin repeat forms a three-helix bundle and was reported primarily in the animal kingdom [30]. These repeats act as modules building long, extended molecules that also serve as a docking surface for cytoskeletal and signal transduction proteins. In *L. pneumophila* it is present in up to eight proteins of each strain (Table 3) and all spectrin repeat proteins are predicted to be secreted Dot/Icm substrates [31-33]. Another interesting domain is the RAS GEF domain that is present in two proteins encoded by strain Paris one of which (Lpp0350) is conserved in the six strains analyzed. Ras-GEFs are small GTPases typically present in eukaryotes that are involved in numerous cellular processes like gene expression, cytoskeleton re-organization, microtubule organization and vesicular and nuclear transport [34]. GEFs (GDP-GTP exchange factors) regulate Rabs, GTP-binding proteins with conserved functions in membrane trafficking [35]. Interestingly, according to the Pfam database Ras-GEF domains in bacteria are only present in *Legionella*, *Parachlamydia acanthamoebae* and *Protochlamydia amoebophila*, all of which are amoeba-associated bacteria.

Coiled-coil domains have been identified previously in the *L. pneumophila* genomes as this motif can be found in all kingdoms of life. However extended coiled-coil domains are largely absent from bacterial genomes but are typical for archaea and eukaryotes. We thus searched the *L. pneumophila* genomes and 29 other genomes of bacterial pathogens or bacteria present in the aquatic environment (Table 4) for proteins with five or more coiled coil domains. Interestingly, *Legionella* spp, *Streptococcus pneumoniae* and *Pseudomonas aeruginosa* contain the highest percentage of proteins with extended coiled-coil domains (6-11 domains) compared to the number of predicted proteins encoded in their genome and only *P. aeruginosa* and *L. pneumophila* encode proteins containing more than 10 coiled-coil domains (Table 4). Most of these *Legionella* proteins are predicted substrates of the Dot/Icm secretion system [31-33,36]. This suggests that large coiled-coil domains are specific adaptations to the eukaryotic cell probably implicated in interactions with host proteins.

c) High selection pressure acts on the Dot/Icm T4SS and its substrates

Central to the pathogenesis of *L. pneumophila* are the *dot/icm* loci, which together direct assembly of a type IV secretion apparatus [37,38]. Although all *L. pneumophila* strains investigated to date contain the complete *dot/icm* loci, sequence variations among the *dot/icm* genes among different *L. pneumophila* strains have been

Table 2 Orthologous eukaryotic like proteins present in the 6 *L. pneumophila* strains and in *L. longbeachae*

Product	Name	<i>L. pneumophila</i> strains										<i>L. lo</i>		
		Paris	Lens	Corby	Lorraine	HL06041035	Philadelphia							
Glucoamylase (Glucan 1,4-alpha-glucosidase)		<i>lpp0489</i>	99.31	<i>lpl0465</i>	98.93	<i>lpc2921</i>	99.38	<i>lpo0482</i>	99.45	<i>lpv0523</i>	99.14	<i>lpg0422*</i>	95.92	<i>llo2801</i>
Putative inosine-uridine nucleoside N-ribohydrolase §		<i>lpp0208</i>	98.81	<i>lpl0206</i>	98.64	<i>lpc0223</i>	98.94							
SidE protein		<i>lpp0304</i>	98.46	<i>lpl0288</i>	95.25	<i>lpc1602</i>	98.35	<i>lpo0273</i>	97.95	<i>lpv0315</i>	98.28	<i>lpg0234*</i>	94.63	
Putative methyltransferase		<i>lpp0358</i>	98.06	<i>lpl0334</i>	89.28	<i>lpc0359</i>	97.93	<i>lpo0334</i>	97.29	<i>lpv0375</i>	97.80	<i>lpg0282</i>	89.28	<i>llo2356</i>
Conserved exported protein of unknown function		<i>lpp0379</i>	99.63	<i>lpl0354</i>	98.53	<i>lpc0380</i>	99.45	<i>lpo0358</i>	99.45		99.82	<i>lpg0301</i>	99.26	
Phosphatidylcholine-hydrolyzing phospholipase C §		<i>lpp0565</i>	99.37	<i>lpl0541</i>	98.66	<i>lpc2843</i>	99.34	<i>lpo0571</i>	99.21	<i>lpv0603</i>	99.29	<i>lpg0502</i>	97.87	<i>llo1329</i>
Phytanoyl-CoA dioxygenase domain-containing protein 1		<i>lpp0578</i>	99.25	<i>lpl0554</i>	98.60	<i>lpc2829</i>	99.46	<i>lpo0586</i>	98.60	<i>lpv0619</i>	99.25	<i>lpg0515*</i>	99.35	<i>llo3224</i>
Leucine-rich repeat protein		<i>lpp1007</i>	97.53			<i>lpc2344</i>	97.87	<i>lpo1029</i>	93.94	<i>lpv1082</i>	97.87	<i>lpg0945*</i>	97.54	
ecto-ATP diphosphohydrolase II	<i>map</i>	<i>lpp1033</i>	98.95	<i>lpl1000</i>	98.78	<i>lpc2316</i>	98.78	<i>lpo1060</i>	98.69	<i>lpv1110</i>	98.86	<i>lpg0971</i>	98.43	<i>llo1247</i>
Major acid phosphatase Map §		<i>lpp1120</i>	99.06	<i>lpl1124</i>	97.92	<i>lpc0577</i>	98.12	<i>lpo1121</i>	97.93	<i>lpv1267</i>	99.06	<i>lpg1119</i>	98.59	<i>llo1016</i>
Pyruvate decarboxylase		<i>lpp1157</i>	99.70	<i>lpl1162</i>	98.87	<i>lpc0618</i>	99.70	<i>lpo1168</i>	98.69	<i>lpv1308</i>	99.70	<i>lpg1155</i>	98.51	
SAM-dependent methyltransferase §		<i>lpp1192</i>	98.38	<i>lpl1198</i>	98.97	<i>lpc0657</i>	99.15	<i>lpo1205</i>	99.06	<i>lpv1346</i>	97.61	<i>lpg1190</i>	99.15	<i>llo1296</i>
Putative 2OG-Fe(II) oxygenase superfamily protein §		<i>lpp1405</i>	100					<i>lpo1449</i>	96.84	<i>lpv1569</i>	100,00	<i>lpg1450</i>	93.74	
Phospholipase C §		<i>lpp1411</i>	100	<i>lpl1573</i>	93.12	<i>lpc0870</i>	100,00	<i>lpo1455</i>	97.61	<i>lpv1576</i>	100,00	<i>lpg1455</i>	97.93	<i>llo1329</i>
Putative mitogen-activated protein kinase	<i>thi</i>	<i>lpp1439</i>	99.12	<i>lpl1545</i>	98.11	<i>lpc0898</i>	97.61	<i>lpo1483</i>	98.93	<i>lpv1609</i>	99.12	<i>lpg1483*</i>	97.67	<i>llo1682</i>
Thiamine biosynthesis protein NMT-1		<i>lpp1522</i>	99.04	<i>lpl1461</i>	97.98	<i>lpc0988</i>	99.04	<i>lpo1583</i>	97.88	<i>lpv1700</i>	99.04	<i>lpg1565</i>	97.47	<i>llo0920</i>
Leucine-rich repeat-containing protein	<i>purC</i>	<i>lpp1567</i>	97.68			<i>lpc1028</i>	98.44			<i>lpv1852</i>	98.91	<i>lpg1602*</i>	97.98	
Phosphoribosylamidoimidazole-succinocarboxamide synthase	<i>mvaB</i>	<i>lpp1647</i>	100	<i>lpl1640</i>	97.76	<i>lpc1106</i>	99.08	<i>lpo1715</i>	98.98	<i>lpv1936</i>	98.16	<i>lpg1675</i>	97.58	<i>llo3277</i>
Hydroxymethylglutaryl-CoA lyase (HMG-CoA lyase) §		<i>lpp1793</i>	99.34	<i>lpl1794</i>	97.13	<i>lpc1274</i>	98.01	<i>lpo1891</i>	99.01	<i>lpv2102</i>	98.68	<i>lpg1830</i>	99.12	<i>llo0113</i>
Putative apyrase		<i>lpp1880</i>	99.47	<i>lpl1869</i>	98.77	<i>lpc1359</i>	99.74	<i>lpo1975</i>	98.86	<i>lpv2179</i>	99.56	<i>lpg1905</i>	95.34	<i>llo1247</i>
Conserved protein of unknown function		<i>lpp1905</i>												
Leucine-rich repeat-containing protein		<i>lpp1940</i>	94.44					<i>lpo2043</i>	93.7	<i>lpv2255</i>	93.88	<i>lpg1958*</i>	92.56	
ZIP metal transporter family protein §		<i>lpp2018</i>	99.60	<i>lpl2013</i>	99.07	<i>lpc1521</i>	99.47	<i>lpo2138</i>	99.34	<i>lpv2339</i>	99.47	<i>lpg2035</i>	99.07	<i>llo2518</i>
Ankyrin repeat protein		<i>lpp2058</i>	99.2	<i>lpl2048</i>	90.42	<i>lpc1566</i>	98.80	<i>lpo2181</i>	98.05					
Conserved protein of unknown function		<i>lpp2061</i>	99.6	<i>lpl2051</i>	95.95	<i>lpc1569</i>	96.80	<i>lpo2185</i>	95.38					
Sphingosine-1-phosphate lyase I		<i>lpp2128</i>	98.84	<i>lpl2102</i>	98.29	<i>lpc1635</i>	99.06	<i>lpo2245</i>	98.62	<i>lpv2428</i>	98.02	<i>lpg2176*</i>	94.02	
Conserved protein of unknown function		<i>lpp2134</i>	100	<i>lpl2109</i>	98.55	<i>lpc1642</i>	100	<i>lpo2253</i>	99.60	<i>lpv2436</i>	100	<i>lpg2182</i>	96.27	
Conserved protein of unknown function		<i>lpp2419</i>	99.84	<i>lpl2298</i>	99.37	<i>lpc2129</i>	100							
Leucine rich repeat protein		<i>lpp2459</i>	98.98	<i>lpl2316</i>	86.85	<i>lpc2085</i>	90.48	<i>lpo2572</i>	99.43	<i>lpv2704</i>	99.32	<i>lpg2392*</i>	97.28	
Putative unspecific monooxygenase		<i>lpp2468</i>	99.47	<i>lpl2326</i>	99.01	<i>lpc2075</i>	98.88	<i>lpo2586</i>	98.15					
Protein kinase-like		<i>lpp2626</i>	94.88	<i>lpl2481</i>	98.85	<i>lpc1906</i>	95.31	<i>lpo2765</i>	98.70	<i>lpv2900</i>	99.13	<i>lpg2556*</i>	99.13	<i>llo2218</i>
Putative methyltransferase		<i>lpp2747</i>	99.25	<i>lpl2620</i>	99	<i>lpc0443</i>	99.37	<i>lpo2974</i>	97.49	<i>lpv3039</i>	99.62	<i>lpg2693</i>	99.37	<i>llo2356</i>
Phytanoyl-CoA dioxygenase, PhyH		<i>lpp2748</i>	99.76	<i>lpl2621</i>	98.91	<i>lpc0442</i>	99.15	<i>lpo2975</i>	98.67	<i>lpv3040</i>	99.76	<i>lpg2694*</i>	95.44	
Sugar kinase §	<i>hemG</i>	<i>lpp2874</i>	99.38			<i>lpc3108</i>	98.89	<i>lpo3114</i>	98.15	<i>lpv3175</i>	99.14	<i>lpg2821</i>	98.52	
Protoporphyrinogen oxidase §	<i>cysK</i>	<i>lpp2909</i>	98.14	<i>lpl2763</i>	96.65	<i>lpc3136</i>	98.90	<i>lpo3153</i>	98.69	<i>lpv3207</i>	98.83	<i>lpg2851</i>	96.02	<i>llo0133</i>
Cysteine synthase A, O-acetylserine sulfhydrylase A subunit		<i>lpp3022</i>	99.26	<i>lpl2880</i>	98.95	<i>lpc3266</i>	95.99	<i>lpo3279</i>	99.16	<i>lpv3334</i>	99.26	<i>lpg2951</i>	98.52	<i>llo0076</i>
Putative methyltransferases §		<i>lpp3025</i>	98.50	<i>lpl2883</i>	97.06	<i>lpc3269</i>	99.30	<i>lpo3282</i>	97.62	<i>lpv3338</i>	97.54	<i>lpg2954</i>	97.76	<i>llo0074</i>

Table 3 Orthologous proteins with eukaryotic motifs present in the 6 *L. pneumophila* strains and in *L. longbeachae*

Motif	<i>L. pneumophila</i> strains										<i>L. lo</i>		
	Paris		Lens		Philadelphia		Lorraine		HL06041035			Corby	
ANK	<i>lpp0037</i>	96.30	<i>lpl0038</i>	97.40	<i>lpg0038*</i>	97.04	<i>lpo0042</i>	97.89	<i>lvp0043</i>	93.66	<i>lpc0039</i>	97.10	
ANK	<i>lpp0126</i>	98.94	<i>lpl0111</i>	98.48	<i>lpg0112</i>	94.83	<i>lpo0119</i>	98.79	<i>lvp0127</i>	93.03	<i>lpc0131</i>	92.16	<i>llo1394</i>
ANK	<i>lpp0202</i>												
ANK	<i>lpp0356</i>												
ANK	<i>lpp0469</i>	98.94	<i>lpl0445</i>	96.35	<i>lpg0403*</i>	95.53	<i>lpo0463</i>	97.64	<i>lvp0501</i>	98.48	<i>lpc2941</i>	98.67	
ANK	<i>lpp0503</i>	98.37	<i>lpl0479</i>	93.86	<i>lpg0436*</i>	93.31	<i>lpo0501</i>	98.12	<i>lvp0537</i>	98.37	<i>lpc2906</i>	98.37	
ANK	<i>lpp0547</i>	99.50	<i>lpl0523</i>	96.31	<i>lpg0483*</i>	96.82	<i>lpo0551</i>	99.83	<i>lvp0585</i>	98.16	<i>lpc2861</i>	99.16	<i>llo2705</i>
ANK	<i>lpp0750</i>	100.00	<i>lpl0732</i>	97.65	<i>lpg0695*</i>	100.00	<i>lpo0775</i>	99.84	<i>lvp0817</i>	100.00	<i>lpc2599</i>	98.44	
ANK	<i>lpp1100</i>												
ANK + SET	<i>lpp1683</i>	97.68	<i>lpl1682</i>	96.32	<i>lpg1718*</i>	98.41	<i>lpo1757</i>	97.86	<i>lvp1985</i>	96.91	<i>lpc1152</i>	97.25	
ANK	<i>lpp1905</i>												
ANK	<i>lpp2058</i>	99.20	<i>lpl2048</i>	90.42			<i>lpo2181</i>	98.05			<i>lpc1566</i>	98.80	
ANK	<i>lpp2061</i>	99.60	<i>lpl2051</i>	95.95			<i>lpo2185</i>	95.38			<i>lpc1569</i>	96.80	
ANK	<i>lpp2065</i>	99.93	<i>lpl2055</i>	98.56			<i>lpo2189</i>	98.62			<i>lpc1573</i>	98.03	
ANK + Fbox	<i>lpp2082</i>	97.40	<i>lpl2072</i>	98.26	<i>lpg2144*</i>	98.84	<i>lpo2207</i>	99.03	<i>lvp2392</i>	93.99	<i>lpc1593</i>	99.22	
ANK	<i>lpp2166</i>	99.25	<i>lpl2140</i>	99.12	<i>lpg2215*</i>	99.06	<i>lpo2285</i>	97.74	<i>lvp2469</i>	99.18	<i>lpc1680</i>	98.93	
ANK	<i>lpp2248</i>	99.50	<i>lpl2219</i>	99.14	<i>lpg2300*</i>	99.14	<i>lpo2371</i>	99.43	<i>lvp2567</i>	98.93	<i>lpc1765</i>	99.21	<i>llo0584</i>
ANK	<i>lpp2270</i>	99.64	<i>lpl2242</i>	97.97	<i>lpg2322*</i>	98.34	<i>lpo2399</i>	98.08	<i>lvp2591</i>	99.53	<i>lpc1789</i>	99.53	<i>llo0570</i>
ANK	<i>lpp2517</i>	99.60	<i>lpl2370</i>	97.94	<i>lpg2452*</i>	98.19	<i>lpo2642</i>	98.95	<i>lvp2776</i>	99.46	<i>lpc2026</i>	99.46	
ANK	<i>lpp2522</i>	98.76	<i>lpl2375</i>	96.90	<i>lpg2456*</i>	95.75	<i>lpo2647</i>	98.49	<i>lvp2781</i>	98.76	<i>lpc2020</i>	91.27	<i>llo0365</i>
ANK	<i>p1pp0098</i>	96.00					<i>lpop0045</i>	96.00					
ANK			<i>lpl1681</i>	100.00							<i>lpc1151</i>	97.98	
ANK			<i>lpl2058</i>	86.17			<i>lpo2193</i>	95.37	<i>lvp2375</i>	94.96			
ANK			<i>lpl2339</i>	98.64	<i>lpg2416*</i>	91.21	<i>lpo2601</i>	99.00	<i>lvp2736</i>	99.28	<i>lpc2057</i>	98.98	
ANK					<i>lpg0402*</i>	100.00			<i>lvp0500</i>	96.01			
ANK									<i>lvp2258</i>				
ANK			<i>lpl1681</i>	100.00							<i>lpc1151</i>	97.98	
ANK			<i>lpl2344</i>	100.00			<i>lpo2607</i>	97.93					
F-Box	<i>lpp0233</i>	98.58	<i>lpl0234</i>	93.97	<i>lpg0171*</i>	96.81	<i>lpo0202</i>	97.87	<i>lvp0254</i>	98.94			
F-Box	<i>lpp2486</i>												
F-Box					<i>lpg2224*</i>	99.83			<i>lvp2482</i>	79.24			
F-Box									<i>lvp2481</i>				
RAS GEF	<i>lpp0350</i>	94.53	<i>lpl0328</i>	96.32	<i>lpg0276*</i>	97.33	<i>lpo0327</i>	97.64	<i>lvp0368</i>	97.64	<i>lpc0353</i>		<i>llo0327</i>
RAS GEFS									<i>lvp2258</i>				
Sec7	<i>lpp1932</i>	98.41	<i>lpl1919</i>	97.40	<i>lpg1950*</i>	92.16	<i>lpo2033</i>	98.32	<i>lvp2243</i>	98.58	<i>lpc1423</i>	97.57	<i>llo1397</i>
Sel1											<i>lpc0165</i>		
Sel1			<i>lpl1059</i>	100.00	<i>lpg1062</i>	99.61			<i>lvp1209</i>	100.00	<i>lpc2212</i>	99.61	
Sel-1§	<i>lpp0957</i>	98.93	<i>lpl0927</i>	98.67	<i>lpg0896</i>	98.93	<i>lpo0978</i>	98.67	<i>lvp1030</i>	98.67	<i>lpc2397</i>	99.47	<i>llo0844</i>
Sel-1	<i>lpp1174</i>	99.39	<i>lpl1180</i>	98.30	<i>lpg1172</i>	98.32	<i>lpo1187</i>	99.11	<i>lvp1327</i>	99.39	<i>lpc0638</i>	99.05	
Sel-1	<i>lpp1310</i>	97.87	<i>lpl1307</i>	98.40	<i>lpg1356</i>	98.67	<i>lpo1345</i>	99.02	<i>lvp1469</i>	97.87	<i>lpc0770</i>	98.76	<i>llo1443</i>
Sel-1	<i>lpp2174</i>	99.64	<i>lpl2147</i>	98.48	<i>lpg2222*</i>	99.56	<i>lpo2292</i>	99.47	<i>lvp2477</i>	99.47	<i>lpc1689</i>	96.27	
Sel-1	<i>lpp2692</i>	99.25	<i>lpl2564</i>	98.61	<i>lpg2639</i>	98.39	<i>lpo2917</i>	99.28	<i>lvp2979</i>	99.39	<i>lpc0501</i>	98.75	<i>llo2649</i>
Sel-1							<i>lpo3233</i>						
Spectrin	<i>lpp1848§</i>	99.18	<i>lpl1845</i>	98.77	<i>lpg1884*</i>	99.01	<i>lpo1944§</i>	98.93	<i>lvp2158§</i>	99.18	<i>lpc1331</i>	99.18	
Spectrin	<i>lpp2246</i>	99.29	<i>lpl2217</i>	98.75	<i>lpg2298*</i>	99.29	<i>lpo2369</i>	99.29	<i>lvp2565</i>	98.27	<i>lpc1763</i>	98.75	<i>llo1707</i>
Spectrin	<i>lpp1930</i>	95.11			<i>lpg1947*</i>	96.65	<i>lpo2029</i>	97.72					
Spectrin	<i>lpp1309</i>	100.00			<i>lpg1355*</i>	90.59			<i>lvp1468</i>	100.00			
Spectrin §	<i>lpp1002</i>	98.01	<i>lpl0971</i>	91.62	<i>lpg0940*</i>	97.92	<i>lpo1024</i>	98.05	<i>lvp1077</i>	97.87	<i>lpc2349</i>	97.15	
Spectrin §	<i>lpp0471</i>	97.79	<i>lpl0447</i>	97.45	<i>lpg0405*</i>	98.30	<i>lpo0465§</i>	98.28	<i>lvp0504</i>	98.28	<i>lpc2939</i>	97.70	<i>llo2845</i>
Spectrin §	<i>lpp1843</i>	95.45	<i>lpl1840</i>	97.57					<i>lvp2151</i>	100.00	<i>lpc1323§</i>	99.60	

Table 3 Orthologous proteins with eukaryotic motifs present in the 6 *L. pneumophila* strains and in *L. longbeachae* (Continued)

Spectrin §	<i>lpp1173</i>	98.56	<i>lpl1179§</i>	98.80	<i>lpg1171*§</i>	98.56	<i>lpo1186</i>	99.28	<i>lpv1326</i>	98.56	<i>lpc0637</i>	97.84	<i>llo3114</i>
STPK	<i>lpp0267</i>	96.95	<i>lpl0262</i>	98.72	<i>lpg0208</i>	93.26	<i>lpo0242</i>	98.92	<i>lpv0288</i>	95.13	<i>lpc0283</i>	97.26	
STPK	<i>lpp1439</i>	99.12	<i>lpl1545</i>	98.11	<i>lpg1483*</i>	97.67	<i>lpo1483</i>	98.93	<i>lpv1609</i>	99.12	<i>lpc0898</i>	97.61	<i>llo1682</i>
STPK	<i>lpp2626</i>	94.88	<i>lpl2481</i>	98.85	<i>lpg2556*</i>	99.13	<i>lpo2765</i>	98.70	<i>lpv2900</i>	99.13	<i>lpc1906</i>	95.31	<i>llo2218</i>
U-box	<i>lpp2887</i>	99.72			<i>lpg2830*</i>	97.15	<i>lpo3124</i>	99.58	<i>lpv3185</i>	98.75			

*Substrates of the Dot/Icm secretion system according to previous publications; § orthologs proteins where the corresponding motif was not present in the other genome; § eukaryotic like proteins newly identified in this study; numbers, nucleotide identity with respect to the *L. pneumophila* Philadelphia gene; *L. lo*, *Legionella longbeachae*

reported [39]. The *dot/icm* loci of the six strains analyzed here exhibited a very high nucleotide conservation of 98-100% among orthologs except for *dotA*, *icmX* and for *icmC* of strain Corby that is shorter and more divergent (84% nucleotide identity) as compared to *icmC* of strain Paris. These results indicate that strong negative selection acts on these genes (Table 5).

Since the identification of RalF [40], numerous approaches have been used to identify Dot/Icm translocated substrates. Currently 278 proteins of *L. pneumophila* have been described as being translocated by the Dot/Icm T4SS system [7,31,32,41-44]. Analysis of their distribution among the six *L. pneumophila* strains reveals a very high conservation, as 206 of the 278 substrates are present in all six strains. Nearly all of them show a nucleotide similarity of 95-100% and only nine are specific to strain Philadelphia (Additional file 1, Table S1). Furthermore, only 34 of the 278 substrates of strain Philadelphia are missing in strain Paris, 30 in strain Lorraine or 25 in strain HL 0604 1035 (Additional file 1, Table S1). Thus, although high redundancy seems to be present in the repertoire of Dot/Icm effectors, the strong conservation of nearly all of them in all genomes, argues for their mutual importance for the *L. pneumophila* life cycle,

Rare exceptions are RalF and AnkB/Lpp2028. The nucleotide sequence of *ralF* of strain Philadelphia is only 85% similar to the *ralF* genes of the other strains and is 72 nts (24aa) shorter. A similar situation is seen for *lpg2144/ankB* that is 54 nts (18aa) longer in strain Philadelphia and Lens than in strain Paris and Corby. This is surprising, as the C-terminal region of AnkB of strain Philadelphia contains a eukaryotic prenylation CAAX motif mediating posttranslational modification of effector proteins, important for intracellular replication of *L. pneumophila*. Lipidation facilitates the localization of this effector protein to host organelles and serves as a docking platform for ubiquitinated proteins [45,46]. Thus in strain Paris and Corby other proteins might take over this function. Taken together, this analysis suggests that over 200 of the Dot/Icm substrates of *L. pneumophila* have been present or have been acquired before the speciation and that such a large repertoire of

effectors is indeed necessary for intracellular replication and adaptation to the specific protozoan hosts.

The species *L. pneumophila* has a highly dynamic accessory genome

a) A wide variety of T4ASSs and conjugative elements contribute to genome plasticity

Based on sequence comparisons, T4SSs are categorized according to their similarity to the *A. tumefaciens* VirB/D4 system into type IVA (type F and P) and type IVB secretion systems [47]. T4ASSs resemble the VirB/D4 system of *A. tumefaciens*, whereas T4BSS proteins are more distantly related to the VirB/D4 proteins [48]. T4SSs are involved in effector translocation, horizontal DNA transfer to other bacteria and eukaryotic cells, in DNA uptake from or release into the extracellular milieu or in the spread of conjugative plasmids [49]. Genome sequence analyses suggest that for *L. pneumophila* T4SSs play an important role for adaptation and virulence as each genome encodes several T4ASSs in addition to the essential T4BSS Dot/Icm discussed above. We identified in each strain either F-type or P-type T4ASSs or both. Figures 2 and Figure 3 show the organization of the structural genes encoding these systems, their organization and their localization (chromosomal or plasmid). The F-type T4ASSs are all predicted to encode a complete T4SS core as well as the essential gene products for pilus assembly and mating pair stabilization that appears to be involved in DNA transfer. They show homology and colinearity with the *tra*-region of the *E. coli* F plasmid [50] and with the recently described *tra* region of *Rickettsia belii* [51]. In *L. pneumophila* strain Philadelphia (Tra5) and *L. longbeachae* strain NSW (Tra6), where the system has a chromosomal localization, it is inserted in a tRNA gene and flanking repeats are present as well as a gene coding for an integrase, suggesting that these T4SSs are mobile (Figure 2). Furthermore, comparison of amino acid identities revealed that the Tra- region on the *L. pneumophila* strain Paris plasmid (Tra1) shows much higher identity with the Tra region located on the *L. longbeachae* plasmid (Tra4) than with those of the different *L. pneumophila* strains (Paris-Tra1, Lens-Tra3 or Lorraine-Tra2)

Table 4 Genes coding for proteins with more than 5 coiled coil domains/protein in different bacterial genomes

Organism	Coiled coil domains proteins	Gene	Product	Number of Coiled coil
<i>B. henselae</i> Houston-1	0			
<i>Ch. pneumoniae</i> J138	0			
<i>Ch. trachomatis</i> D UW-3	0			
<i>C. glutamicum</i> ATCC 13032	0			
<i>E. coli</i> O157:H7	1	<i>ECH74115_2173</i>	tail length tape measure protein	5
<i>H. influenzae</i> Rd KW20	0			
<i>H. pylori</i> 26695	1	<i>HP0527</i>	cag pathogenicity island protein Y	10
<i>L. pneumophila</i> Corby	7	<i>lpc1130</i>	substrate of the Dot/Icm system/lcm system	5
		<i>lpc1131</i>	substrate of the Dot/Icm system/lcm system	6
		<i>lpc1452</i>	substrate of the Dot/Icm system/lcm system	6
		<i>lpc1611</i>	hypothetical protein	12
		<i>lpc1987</i>	substrate of the Dot/Icm system, effector protein B	9
		<i>lpc2349</i>	substrate of the Dot/Icm system, LidA	6
		<i>lpc3079</i>	substrate of the Dot/Icm system, effector protein A	5
<i>L. pneumophila</i> HL06041035	10	<i>lpv1077</i>	substrate of the Dot/Icm system, LidA	6
		<i>lpv1725</i>	substrate of the Dot/Icm system	6
		<i>lpv1966</i>	substrate of the Dot/Icm system	5
		<i>lpv1967</i>	substrate of the Dot/Icm system	6
		<i>lpv2269</i>	substrate of the Dot/Icm system	7
		<i>lpv2408</i>	conserved protein of unknown function	5
		<i>lpv2816</i>	substrate of the Dot/Icm system, effector protein B	10
		<i>lpv2959</i>	chromosome segregation SMC protein	9
		<i>lpv3144</i>	substrate of the Dot/Icm system, effector protein A	5
		<i>lpv3184</i>	substrate of the Dot/Icm system, SidH	9
<i>L. pneumophila</i> Lens	7	<i>lpl1437</i>	substrate of the Dot/Icm system	6
		<i>lpl1660</i>	substrate of the Dot/Icm system	7
		<i>lpl1661</i>	substrate of the Dot/Icm system	6
		<i>lpl1941</i>	substrate of the Dot/Icm system	5
		<i>lpl2084</i>	substrate of the Dot/Icm system	5
		<i>lpl2411</i>	substrate of the Dot/Icm system, effector protein B	9
		<i>lpl2708</i>	substrate of the Dot/Icm system, effector protein A	5
<i>L. pneumophila</i> Lorraine	10	<i>lpo1024</i>	substrate of the Dot/Icm system, LidA	6
		<i>lpo1608</i>	substrate of the Dot/Icm system	6
		<i>lpo1735</i>	substrate of the Dot/Icm system	7
		<i>lpo1736</i>	substrate of the Dot/Icm system	5
		<i>lpo2060</i>	substrate of the Dot/Icm system	6
		<i>lpo2216</i>	substrate of the Dot/Icm system, SdeC	5
		<i>lpo2680</i>	substrate of the Dot/Icm system, effector protein B	9
		<i>lpo2896</i>	chromosome segregation SMC protein	9
		<i>lpo3083</i>	substrate of the Dot/Icm system, effector protein A	5
<i>lpo3123</i>	substrate of the Dot/Icm system	9		
<i>L. pneumophila</i> Paris	6	<i>lpp1002</i>	substrate of the Dot/Icm system, LidA	6
		<i>lpp1546</i>	substrate of the Dot/Icm system	6

Table 4 Genes coding for proteins with more than 5 coiled coil domains/protein in different bacterial genomes (Continued)

		<i>lpp1666</i>	substrate of the Dot/Icm system	7
		<i>lpp1952</i>	substrate of the Dot/Icm system	6
		<i>lpp2555</i>	substrate of the Dot/Icm system, effector protein B	10
		<i>lpp2883</i>	substrate of the Dot/Icm system	6
<i>L. pneumophila</i> Philadelphia	8	<i>lpg1355</i>	substrate of the Dot/Icm system, SidG protein	5
		<i>lpg1588</i>	substrate of the Dot/Icm system	6
		<i>lpg1701</i>	substrate of the Dot/Icm system	5
		<i>lpg1702</i>	substrate of the Dot/Icm system	6
		<i>lpg2156</i>	protein of unknown function	5
		<i>lpg2490</i>	substrate of the Dot/Icm system, effector protein B	9
		<i>lpg2793</i>	substrate of the Dot/Icm system, effector protein A	5
		<i>lpg2829</i>	substrate of the Dot/Icm system	8
<i>L. monocytogenes</i> EGD-e	3	<i>lmo0650</i>	hypothetical protein	5
		<i>lmo0955</i>	hypothetical protein	5
		<i>lmo1224</i>	hypothetical protein	5
<i>M. tuberculosis</i> F11	1	<i>TBFG_12936</i>	chromosome partitioning protein Smc	10
<i>M. tuberculosis</i> H37Ra	1	<i>MRA_2947</i>	putative chromosome segregation Smc	10
<i>N. meningitidis</i> MC58	0			
<i>P. aeruginosa</i> LESB58	11	<i>PLES_08211</i>	putative tail length tape measure protein	7
		<i>PLES_12531</i>	hypothetical protein	7
		<i>PLES_12541</i>	hypothetical protein	5
		<i>PLES_13581</i>	putative tail length tape measure protein	7
		<i>PLES_15241</i>	electron transport complex protein RnfC	8
		<i>PLES_15871</i>	hypothetical protein	6
		<i>PLES_36651</i>	putative ClpA	—
		<i>PLES_38011</i>	putative chromosome segregation protein	11
		<i>PLES_46621</i>	putative exonuclease	13
		<i>PLES_50721</i>	hypothetical protein	6
		<i>PLES_55491</i>	putative outer membrane protein precursor	5
<i>R. felis</i> URRWXCal2	2	<i>RF_0022</i>	putative surface cell antigen sca1	7
		<i>RF_0725</i>	antigenic heat-stable 120 kDa protein	5
<i>R. prowazekii</i> Madrid E	0			
<i>R.a typhi</i> Wilmington	0			
<i>S. typhimurium</i> LT2	5	<i>STM0395</i>	exonuclease subunit SbcC	7
		<i>STM0567</i>	putative DNA repair ATPase	7
		<i>STM0994</i>	chromosome partition protein mukB	10
		<i>STM1041</i>	minor tail protein	5
		<i>STM3199</i>	hypothetical protein	5
<i>S. flexneri</i> 2a 2457T	1	<i>S0984</i>	fused chromosome partitioning protein	10
<i>Synechocystis</i> sp. PCC 6803	2	<i>sl11772</i>	MutS2 protein	5
		<i>slr1301</i>	hypothetical protein	6
<i>S. pneumoniae</i> D39	4	<i>SPD_0126</i>	exported protein of unknown function	6
		<i>SPD_0710</i>	putative Septation ring formation regulator EzrA	7
		<i>SPD_1104</i>	chromosome partition protein Smc	10
		<i>SPD_2017</i>	exported protein of unknown function	6
<i>W. pipientis</i> wMel	0			

Table 4 Genes coding for proteins with more than 5 coiled coil domains/protein in different bacterial genomes (Continued)

<i>X. fastidiosa</i> 9a5c	0			
<i>Y. pestis</i> KIM	4	y0227	hypothetical protein	6
		y0976	ATP-dependent dsDNA exonuclease	12
		y2765	chromosome partition protein MukB	10
		yapB	autotransporter	6

(Figure 2). Thus these systems seem to be transferred horizontally via plasmids but are also able to integrate in the genome similar to what was reported for the Lvh-region [52].

The F-type T4SS encode long, flexible pili that allow donors to mate in liquid and on solid media with equal efficiencies [53]. In contrast P-type T4SS like described in *P. aeuroginosa* encode short and rigid conjugative pili that allow surface mating. Homologues to this system are also present in the *Legionella* genomes. They were initially described in two genomic islands of *L. pneumophila* strain Corby (Figure 3; Trb1 and Trb2) [54]. We show here that they are also present in the chromosomes of *L. pneumophila* strain Lorraine (Trb3) and *L.*

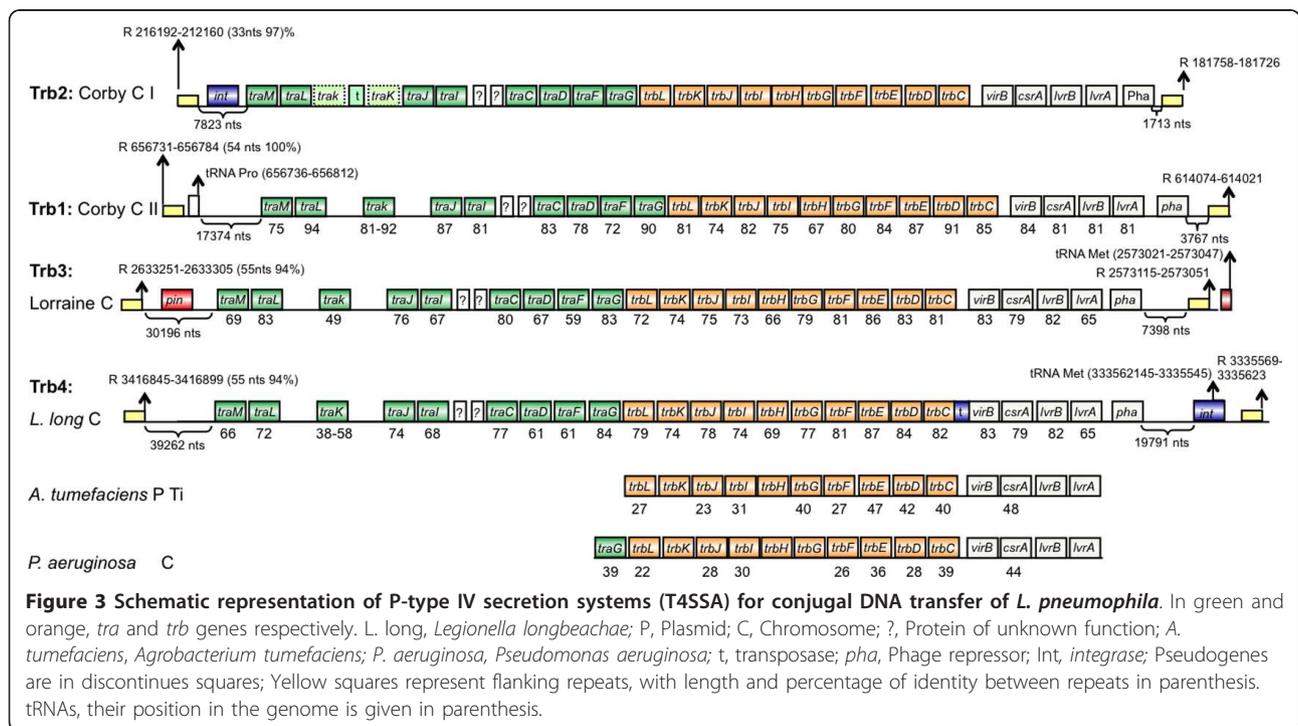
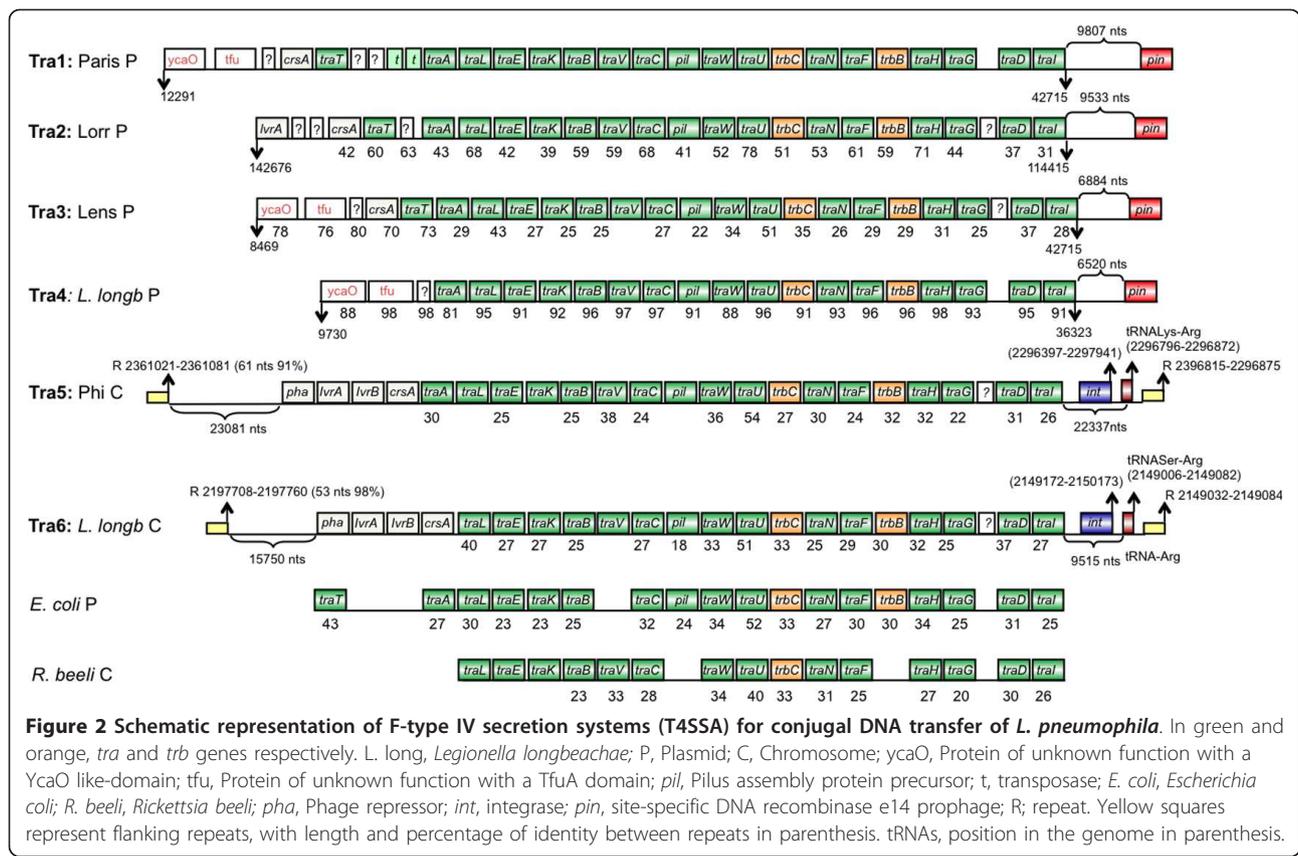
longbeachae NSW150 (Trb4) (Figure 3). Again for all T4SS regions flanking repeats are found suggesting mobility, and protein identity values and GC-content values of the *tra-trb* genes are higher than the genomic average (38%), supporting again horizontal and not vertical transmission.

Another intriguing feature of these regions is that several transposases and phage related proteins are present in each of the *tra* clusters as well as genes coding for homologues of a putative phage repressor protein (PrpA) and for homologues of LvrA, LvrB and LvrC, first described for the Lvh region of *L. pneumophila*. LvrC is a homologue of CsrA, a protein crucial for the regulation of the switch between replicative and

Table 5 Percentage of nucleotide identity of orthologous *dot/icm* genes with respect to the *L. pneumophila* Philadelphia sequence

Gene name	Length (nts)	Phila	Paris	Id	Lens	Id	Lorrain	Id	HL06041035	Id	Corby	Id	<i>L. long</i>	Id
<i>icmT</i>	261	<i>lpg0441</i>	<i>lpp0507</i>	99.6	<i>lpl0483</i>	99.1	<i>lpo0507</i>	100	<i>lpv0541</i>	96	<i>lpc2902</i>	99.2	<i>llo2795</i>	75.2
<i>icmS</i>	345	<i>lpg0442</i>	<i>lpp0508</i>	98.5	<i>lpl0484</i>	98.8	<i>lpo0508</i>	99.1	<i>lpv0542</i>	94.4	<i>lpc2901</i>	98.3	<i>llo2794</i>	76.9
<i>icmR</i>	363	<i>lpg0443</i>	<i>lpp0509</i>	96.9	<i>lpl0485</i>	98.3	<i>lpo0509</i>	97.8	<i>lpv0543</i>	97.5	<i>lpc2900</i>	96.9		
<i>icmQ</i>	576	<i>lpg0444</i>	<i>lpp0510</i>	97	<i>lpl0486</i>	99	<i>lpo0510</i>	98	<i>lpv0544</i>	98	<i>lpc2899</i>	98	<i>llo2792</i>	70.7
<i>icmP/dotM</i>	1131	<i>lpg0445</i>	<i>lpp0511</i>	98	<i>lpl0487</i>	99	<i>lpo0511</i>	98	<i>lpv0545</i>	98	<i>lpc2898</i>	99	<i>llo2791</i>	74.5
<i>icmO/dotL</i>	2352	<i>lpg0446</i>	<i>lpp0512</i>	98.4	<i>lpl0488</i>	97.7	<i>lpo0512</i>	98.1	<i>lpv0546</i>	98.3	<i>lpc2897</i>	98.3	<i>llo2790</i>	77.7
<i>icmN/DotK</i>	570	<i>lpg0447</i>	<i>lpp0513</i>	99.3	<i>lpl0489</i>	98.6	<i>lpo0513</i>	98.9	<i>lpv0547</i>	99.6	<i>lpc2896</i>	99.7	<i>llo2789</i>	67.3
<i>icmM/dotJ</i>	285	<i>lpg0448</i>	<i>lpp0514</i>	97.9	<i>lpl0490</i>	97.9	<i>lpo0514</i>	97.9	<i>lpv0548</i>	99.3	<i>lpc2895</i>	98.6	<i>llo2788</i>	61.7
<i>icmL/dotI</i>	639	<i>lpg0449</i>	<i>lpp0515</i>	99.8	<i>lpl0491</i>	99.4	<i>lpo0515</i>	99.4	<i>lpv0549</i>	99.8	<i>lpc2894</i>	99.5	<i>llo2787</i>	78.6
<i>icmK/dotH</i>	1083	<i>lpg0450</i>	<i>lpp0516</i>	94.8	<i>lpl0492</i>	94.3	<i>lpo0516</i>	95.2	<i>lpv0550</i>	94.4	<i>lpc2893</i>	94.7	<i>llo2786</i>	71.2
<i>icmE/dotG</i>	3147	<i>lpg0451</i>	<i>lpp0517</i>	93.7	<i>lpl0493</i>	94.0	<i>lpo0517</i>	94	<i>lpv0551</i>	94	<i>lpc2892</i>	94.3	<i>llo2785</i>	69.1
<i>icmG/dotF</i>	810	<i>lpg0452</i>	<i>lpp0518</i>	98	<i>lpl0494</i>	97	<i>lpo0518</i>	98	<i>lpv0552</i>	98	<i>lpc2891</i>	97	<i>llo2784</i>	55.7
<i>icmC/dotE</i>	585	<i>lpg0453</i>	<i>lpp0519</i>	99.6	<i>lpl0495</i>	99.1	<i>lpo0519</i>	99.7	<i>lpv0553</i>	99.3	<i>lpc2890</i>	54	<i>llo2783</i>	69.1
<i>icmD/DotP</i>	399	<i>lpg0454</i>	<i>lpp0520</i>	97	<i>lpl0496</i>	98	<i>lpo0520</i>	97	<i>lpv0554</i>	98	<i>lpc2889</i>	97	<i>llo2782</i>	77.3
<i>icmJ/dotN</i>	627	<i>lpg0455</i>	<i>lpp0521</i>	99	<i>lpl0497</i>	98	<i>lpo0521</i>	99	<i>lpv0555</i>	99	<i>lpc2888</i>	98	<i>llo2781</i>	79.4
<i>icmB/DotO</i>	3030	<i>lpg0456</i>	<i>lpp0522</i>	98.1	<i>lpl0498</i>	98.3	<i>lpo0522</i>	98.3	<i>lpv0556</i>	98.2	<i>lpc2887</i>	97.6	<i>llo2780</i>	76.4
<i>icmF</i>	2922	<i>lpg0458</i>	<i>lpp0524</i>	98.2	<i>lpl0500</i>	98.5	<i>lpo0524</i>	98.3	<i>lpv0558</i>	98.5	<i>lpc2885</i>	98.2	<i>llo3075</i>	69.5
<i>icmH/DotU</i>	786	<i>lpg0459</i>	<i>lpp0525</i>	99.4	<i>lpl0501</i>	99.5	<i>lpo0525</i>	99.7	<i>lpv0559</i>	99	<i>lpc2884</i>	99	<i>llo3074</i>	68.8
<i>dotD</i>	492	<i>lpg2674</i>	<i>lpp2728</i>	98	<i>lpl2601</i>	98	<i>lpo2953</i>	98	<i>lpv3018</i>	98	<i>lpc0463</i>	99	<i>llo0369</i>	76.5
<i>dotC</i>	912	<i>lpg2675</i>	<i>lpp2729</i>	98.7	<i>lpl2602</i>	98.5	<i>lpo2954</i>	98.8	<i>lpv3019</i>	98.6	<i>lpc0462</i>	99.9	<i>llo0368</i>	74.8
<i>dotB</i>	1134	<i>lpg2676</i>	<i>lpp2730</i>	99	<i>lpl2603</i>	98	<i>lpo2955</i>	98	<i>lpv3020</i>	98	<i>lpc0461</i>	99	<i>llo0367</i>	76
<i>dotA</i>	3108	<i>lpg2686</i>	<i>lpp2740</i>	83.3	<i>lpl2613</i>	96.8	<i>lpo2967</i>	83	<i>lpv3032</i>	83.6	<i>lpc0450</i>	85.8	<i>llo0364</i>	51.4
<i>icmV</i>	456	<i>lpg2687</i>	<i>lpp2741</i>	91	<i>lpl2614</i>	91	<i>lpo2968</i>	91	<i>lpv3033</i>	92	<i>lpc0449</i>	92	<i>llo0363</i>	64.3
<i>icmW</i>	456	<i>lpg2688</i>	<i>lpp2742</i>	95.1	<i>lpl2615</i>	97.6	<i>lpo2969</i>	95.1	<i>lpv3034</i>	95.4	<i>lpc0448</i>	95.1	<i>llo0362</i>	79.3
<i>icmX</i>	1419	<i>lpg2689</i>	<i>lpp2743</i>	84.3	<i>lpl2616</i>	85.2	<i>lpo2970</i>	85.6	<i>lpv3035</i>	85.6	<i>lpc0447</i>	84.1	<i>llo0361</i>	46.9

Id, identity



transmissible phase of *L. pneumophila* [55]. It is tempting to assume that these CsrA homologues are implicated in the regulation of the mobility of these islands. Possibly, dependent on the growth phase and/or on metabolic cues *L. pneumophila* might excise these islands as multiple copies could be advantageous in certain conditions, or perhaps allow high frequencies of DNA transfer leading to fast and efficient adaptation to new conditions. The genomic features of these islands suggest a particular mechanism of mobility, which will be interesting to investigate.

b) The *L. pneumophila* genomes encode systems specific for protection against invading DNA and stabilization of large genomic fragments

Bacteria have developed multiple methods of protection against mobile genetic elements or bacteriophages. An example for acquired phage specific immunity is clustered regularly interspaced short palindromic repeats (CRISPR) loci [56]. Another type of protection may be conferred by toxin-antitoxin (TA) systems. Bacterial TA systems are small genetic modules composed of a toxin and antitoxin. While toxins are always proteins,

antitoxins are either RNAs (type I and III) or proteins (type II) [57]. These systems were first described for being dedicated to plasmid maintenance. Several lines of research indicate that chromosomal TA systems might serve as protection against mobile genetic elements such as plasmids and phages. However, recent studies have shown that type II systems are also involved in the stabilization of large genomic fragments and of integrative conjugative elements [57]. Interestingly, type II TA systems are thought itself to be part of the mobilome and to move from one genome to another through horizontal gene transfer [57].

Genome analyses identified several TA and CRISPR systems. Interestingly, we identified only type II TA systems of which all except two are in a chromosomal location (Table 6). However, of the 18 chromosomal encoded TA systems identified at least 14 are located on putative genomic islands or mobile genetic elements. The two most frequently found TA systems in the *L. pneumophila* genomes are homologues of the HigAB and RelEB systems. HigAB was first described in the *Vibrio cholerae* superintegron where it encodes mRNA

Table 6 Genes encoding putative toxin-antitoxin systems in six *L. pneumophila* genomes

Toxin-antitoxin	<i>L. pneumophila</i> strains					
	Paris	Lens	Philadelphia	Corby	Lorraine	HL06041035
<i>higA</i>		<i>lpl2833</i> (96)*	<i>lpg2914</i> (96)			<i>lpv3285</i> (96)
<i>higB</i>		<i>lpl2834</i> (87)*	<i>lpg2915</i> (103)			<i>lpv3286</i> (103)
<i>higA</i>		<i>lpl1092</i> (93)*				
<i>higB</i>		<i>lpl1093</i> (107)*				
<i>higA</i>	<i>lpp0064</i> (434)*				<i>lpo0072</i> (432)*	
<i>higB</i>	<i>lpp0065</i> (79)*				<i>lpo0073</i> (79)*	
				<i>lpc2112</i> (312)		
		<i>lpl2291</i> (102)*	<i>lpg2369</i> (102)	<i>lpc2113</i> (37)		<i>lpv2676</i> (102)*
Similar to <i>hipA</i>	<i>lpp2427</i> (78)*	<i>lpl2292</i> (312)*	<i>lpg2370</i> (312)	<i>lpc2114</i> (65)	<i>lpo2551</i> (115)*	<i>lpv2677</i> (310)*
<i>yhvA</i>					<i>lpo1074</i> (168)*	
<i>sohA</i>					<i>lpo1075</i> (115)*	
<i>relE</i>	<i>p1pp0090</i> (83)	<i>lpl1587</i> (82)*				
<i>relB</i>	<i>p1pp0089</i> (95)	<i>lpl1588</i> (85)*				
<i>relE</i>		<i>lpl1084</i> (84)*		<i>lpc2177</i> (93)*	<i>lpo0120</i> (93)*	
<i>relB</i>				<i>lpc2178</i> (88)*	<i>lpo0119</i> (86)*	
<i>parE</i>						<i>lpe2361</i> (98)*
<i>parD</i>						<i>lpe2360</i> (84)*
<i>pemK</i>					<i>lpo0114</i> (106)	

*TA systems located on putative genomic islands; in parenthesis length of the corresponding protein

cleaving enzymes and can stabilize plasmids [58]. RelEB was shown, when introduced into the *E. coli* chromosome to prevent deletion of flanking DNA and thus to diminish large scale genome reduction [59]. The same function was shown for the ParED system of *Vibrio vulnificus*, homologues of which are also present in one of the *L. pneumophila* genomes (Table 6). Thus, the different *L. pneumophila* TA systems might be important for stabilization of plasmids and integrative conjugative elements and for protection against invasion of plasmids, phages, or other mobile genetic elements.

The CRISPR/cas system was shown to provide resistance against invading viruses and plasmids and has been identified in many bacteria and archaea [60]. CRISPR/cas loci are also present in the *L. pneumophila* genomes of strains Paris, Lens, Alcoy and 130 b but are absent from strains HL06041035 and Lorraine. According to the cas genes, the CRISPR locus of Paris is closely related to that of strain 130 b. In contrast the one of strain Lens located on the plasmid is closely related to the chromosomal CRISPR locus of strain Alcoy as previously described [61]. Strain Lens carries a second CRISPR locus on the chromosome; however, it does not seem to be functional like the one encoded by strain Alcoy. Probably strong protection against invading phages is not extremely important, as not all *L. pneumophila* strains contain CRISPR loci. This may be related to their intracellular life style or that despite their widespread occurrence in aquatic environments only few bacteriophages that specifically infect *Legionella* seem to exist [62].

c) Accessory genome of strains Lorraine and HL 0604 1035

In order to get insight in the genetic basis of the two newly sequenced strains, possibly implicated in their different disease frequencies (Lorraine is a newly emerging endemic clone and strain HL 0604 1035 is a *L. pneumophila* Sg1 strain never isolated from disease) we analyzed the specific gene content of each of these strains more in depth. Strain HL 0604 1035 contains 92 and strain Lorraine 148 genes without homology to any gene of the other five *L. pneumophila* strains sequenced of which the majority (60 in strain HL 0604 1035 and 73 in strain Lorraine) code for proteins of unknown function (Additional file 2, Tables S2 and additional file 3, Table S3). Among the genes in these two genomes that lack an ortholog in the other sequenced *L. pneumophila* genomes, about 50% are clustered on three large genomic islands. One genomic Island (GI-HL1) of 45 kb spans from *lpv2637* to *lpv2691*. It is bordered by a Met tRNA gene and encodes a phage related integrase. A second putative mobile element (GI-HL2) of 27 kbs contains the region from *lpv0193* to *lpv0226*. It is bordered at one side by an integrase and a reverse transcriptase (*lpv0225*) and on the other side by a prophage

Rac integrase and a phage excisionase. Strain Lorraine contains also a large genomic island (GI-Lo1) of 69 kb that spans from *lpo2442* to *lpo2531*. It is inserted in a Met tRNA gene, contains a phage related integrase and flanking repeats of 72 nts. Additional, smaller genomic islands seem to be present, however, their borders are difficult to define. Thus most of the strain specific genes seem to be acquired by HGT through mobility of genomic islands.

Only for few of the specific genes a putative function can be predicted like genes coding for proteins involved in sugar and nucleotide metabolism, for uridine diphosphoglucuronate 5'-epimerase or for an UDP-glucose 6-dehydrogenase. Furthermore a specific ANK motif containing protein and a leucine rich repeat protein are present in strain HL 0604 1035. In strain Lorraine we identified mainly specific metabolic enzymes like a putative flavanone 3-dioxygenase, an enzyme involved in flavonoids metabolism and in biosynthesis of phenylpropanoids, which are secondary metabolites of plants and algae. In addition, *lpo2614* is predicted to encode a kynurenine-oxoglutarate transaminase, an enzyme that is part of the tryptophan metabolism and *lpo2960* codes for a putative glycolate oxidase that catalyses the conversion of glycolate and oxygen to glyoxylate and hydrogen peroxide. *lpo2502* codes a homologue of CsbD, a general stress response protein of *Bacillus subtilis* [63]. However, the best BLASTp hit is with the *Protochlamydia amoebophila* homologue, an *Acanthamoeba* sp. symbiont [64]. Probably this gene has been acquired by HGT between these two bacteria within their amoeba host. Quite surprisingly, we identified a gene coding a putative methyl-accepting chemotaxis sensory transducer (*lpv1770*) although all *L. pneumophila* strains analyzed to date do not encode chemotaxis systems. This gene shares 71.34% amino acid identity with Llo3301 of *L. longbeachae* a protein that is part of its chemotaxis system [28] also present in *L. drancourtii* [65]. Probably a common ancestor encoded a chemotaxis system that was lost in *L. pneumophila* through a deletion and degradation process.

d) Shared genome of the epidemic strains Paris and Lorraine

A search for genes shared by the two endemic strains but absent in all other strains identified only three genes that fulfilled these criteria and for which a function could be predicted. These encode the alpha, beta and gamma subunits of a putative thiocyanate hydrolase (*lpo1236*, *lpo1237*, *lpo1238* and *lpp1219*, *lpp1220*, *lpp1221*). Most interestingly, these strains are both common in France and strain Paris is also world-wide distributed [10] suggesting a better niche adaptation. Indeed, thiocyanate compounds are used for cleaning water circuits and these strains are thus probably able to

better resist these treatments [66]. Furthermore, strain Alcoy that is responsible for several outbreaks and many cases of Legionnaires' disease in Spain, also contains these genes [61]. The genes coding the putative thiocyanate hydrolase have a GC content of 41-43%, which is significantly higher than the average G+C content of the *L. pneumophila* genome, which is 38%. When searching for the closest homologues according to BLAST searches we identified them in the genomes of *Rhodococcus opacus* strain B4 and *Nocardia farcinica* spp. These two are high G+C Gram-positive bacteria belonging to the *Actinomycetales*, which are phylogenetically not closely related to *Legionella* suggesting that *L. pneumophila* acquired these genes by horizontal gene transfer.

Taken together, the analysis of the accessory gene content showed again that *L. pneumophila* genomes show high plasticity due to mobile genetic elements and

HGT. No specific virulence related genes explaining their different disease frequencies have been identified. However, the identification of a specific thiocyanate hydrolase might explain the wide distribution of strains Paris and Lorraine as it may allow them to better adapted to artificial water systems.

Evolutionary genomics

Phylogenetic reconstruction reveals extensive recombination

To analyze the relationship among the six different *L. pneumophila* strains a phylogenetic reconstruction was done based on a multilocus sequence (MLSA) approach using 31 genes selected according to Zeigler [67] (Table 7 and Additional file 4, Table S4). These 31 genes were chosen as they had been shown to be powerful for predicting the relatedness of bacterial genomes [67]. The phylogeny obtained from their concatenated alignment showed a well-resolved topology with bootstrap values

Table 7 Characteristics of the 31 genes used for phylogenetic reconstruction

Gene Name	Product	Label ^a	Function	Length (nts) ^a
<i>uvrB</i>	Excinuclease ABC, subunit B	<i>lpp0086</i>	DNA replication, recombination, and repair	1992
<i>pgk</i>	Phosphoglycerate kinase	<i>lpp0152</i>	Glycolysis/gluconeogenesis	1191
<i>rpoA</i>	RNA polymerase, alpha subunit	<i>lpp0419</i>	Transcription	993
<i>ffh</i>	Signal recognition particle protein, GTPase	<i>lpp0467</i>	Transport and binding proteins	1377
<i>serS</i>	Seryl tRNA synthetase	<i>lpp0575</i>	tRNA aminoacylation	1281
<i>proS</i>	Prolyl-tRNA synthase	<i>lpp0749</i>	tRNA aminoacylation	1710
<i>glyA</i>	Serine hydroxymethyltransferase	<i>lpp0791</i>	Glycine/serine hydroxymethyltransferase	1254
<i>dnaB</i>	Replicative DNA helicase	<i>lpp0803</i>	DNA replication, recombination, and repair	1383
<i>gpi</i>	Glucose-6-phosphate isomerase	<i>lpp0825</i>	Glycolysis/gluconeogenesis	1500
<i>lig</i>	DNA ligase	<i>lpp1020</i>	DNA replication, recombination, and repair	2022
<i>cysS</i>	Cysteinylyl-tRNA synthetase	<i>lpp1271</i>	tRNA aminoacylation	1371
<i>trpS</i>	Tryptophanyl tRNA synthetase	<i>lpp1399</i>	tRNA aminoacylation	1215
<i>aspS</i>	Aspartyl-tRNA synthetase	<i>lpp1434</i>	tRNA aminoacylation	1782
<i>ruvB</i>	Holliday junction DNA helicase	<i>lpp1534</i>	tRNA aminoacylation	1011
<i>nrdA</i>	Ribonucleoside-diphosphate reductase, alpha subunit	<i>lpp1738</i>	Deoxyribonucleotide/ribonucleoside metabolism	2829
<i>recA</i>	Bacterial DNA recombination protein	<i>lpp1765</i>	DNA replication, recombination, and repair	1047
<i>tig</i>	Trigger factor	<i>lpp1830</i>	Protein folding and stabilization	1332
<i>lepA</i>	GTP-binding membrane protein	<i>lpp1837</i>	Translation	1833
<i>metK</i>	S-adenosylmethionine synthetase	<i>lpp2004</i>	tRNA aminoacylation	1149
<i>dnaJ</i>	Heat shock protein	<i>lpp2006</i>	Protein folding and stabilization	1140
<i>argS</i>	Arginyl tRNA synthetase	<i>lpp2013</i>	tRNA aminoacylation	1770
<i>eno</i>	Enolase	<i>lpp2020</i>	Glycolysis/gluconeogenesis	1269
<i>ftsZ</i>	Cell division protein	<i>lpp2662</i>	Cell division	1197
<i>uvrC</i>	Excinuclease ABC, subunit C	<i>lpp2698</i>	DNA replication, recombination, and repair	1857
<i>dnaX</i>	DNA polymerase III, subunits gamma and tau	<i>lpp2802</i>	DNA replication, recombination, and repair	1671
<i>recN</i>	DNA repair protein	<i>lpp2877</i>	DNA replication, recombination, and repair	1668
<i>metG</i>	Methionyl tRNA synthetase	<i>lpp2941</i>	tRNA aminoacylation	2013
<i>rho</i>	Transcription terminator factor	<i>lpp3002</i>	Translation	1262
<i>atpD</i>	ATP synthase F1, subunit beta	<i>lpp3053</i>	ATP-proton motive force interconversion	1377
<i>atpA</i>	ATP synthase, subunit alpha	<i>lpp3055</i>	ATP-proton motive force interconversion	1554
<i>thdF</i>	GTP binding protein, thiophene oxidation	<i>lpp3073</i>	tRNA and rRNA base modification	1341

^a with respect to strain Paris, nts nucleotides

over 50%. To ascertain the reliability of the obtained phylogenetic tree we established individual phylogenies for each of the 31 genes. Surprisingly, the incongruence among several gene trees was high. In addition the Consense program results did not support any node to at least 50%. To further investigate these results we undertook a second analysis using a Shimodaira-Hasegawa test and compared the topologies of the individual alignments of each gene and the concatenated alignment of the 31 genes. As shown in Additional file 5, Table S5 the likelihood-based SH test for alternative tree topologies identified striking discordances. A possible explanation for the identified incongruences among the phylogenies obtained in our study is the presence of recombination events.

With the aim to explore whether recombination events are present in the selected genes we undertook an in depth analysis using the program RDP [68]. Indeed, the analysis of individual genes identified intragenic recombination in 9 of the 31 genes (Table 8). Numerous additional recombination events were detected with the concatenated alignment of the 22 genes for which no intragenic recombination had been shown (Table 8). To minimize false positive recombination events only those that were supported by at least two of the six methods used in RDP were taken into account. However, except one, all were supported by at least three methods. No artifacts resulting of positive selection should be included in this analysis since all of the genes are either informational or operational (housekeeping). Most interestingly, four of the genes in which intragenic recombination was detected are housekeeping genes (*pgk*, *atpD*, *ffh*, *metK*). Housekeeping genes allow to estimate the extent of recombination within bacterial species since presence of recombination in such “normally recombination free genes” is indicative of a high rate of recombination [22]. Similarly antigen-coding genes of *Legionella* were reported to show recombination events [18,69] and certain other genomic regions [17,19,70-72]. Another example of intragenic recombination in *L. pneumophila* is the *rtxA* gene that contains a long tandem repeated domain of variable copy number and sequence [4,10,73]. *rtxA* of strain Lorraine and Corby share the same repeats, whereas the other strains have unique types of repeats. However, when including the newly sequenced strains Lorraine and HL 0604 1035 we found that repeats of the same type are shared by HL 0604 1035 and Philadelphia and by Lorraine and Lens (Figure 4 and Additional file 6, Table S6), further substantiating high intragenic recombination among strains.

To reconstruct the phylogenetic history of the species *L. pneumophila* we used thus the concatenated alignment of the 31 genes described above. It gave a topology with high bootstrap support, however recombination

bias may result in high support for the wrong tree. To avoid possible bias we thus analyzed the concatenated alignment of the 31 genes using a split tree decomposition that allows a more realistic representation of the phylogenetic relationships. Furthermore we constructed a classical bifurcating tree using the highest possible number of genes [all orthologs among the six strains with (1867 genes) and without (2434 genes) *L. longbeachae* as outgroup]. As shown in Figure 5 the Splits Decomposition phylogeny is network-like suggesting incompatible partitions within sequence data, which commonly arise from recombination. Although the phylogeny based on the orthologous genes can also be affected by recombination, the high number of informative sites included in this data set, should allow recovering the correct history of the species as it has been shown previously for other closely related bacterial species [74].

Taken together, in contrast to previous studies, which reported that the species *L. pneumophila* is a clonal population [13,14] our results show clearly that a high recombination rate shapes the *L. pneumophila* genomes. This finding is in line with the natural competence of *L. pneumophila*. However, some worldwide distributed *L. pneumophila* clones have been described (e.g. [10]), suggesting that *L. pneumophila* is able to develop a unique genetic population structure within a particular region or environment as reported recently [72].

Recombination of large chromosomal regions of over 200 kbs among *L. pneumophila* strains

Our recombination analysis revealed not only intragenic recombination events but also intergenic recombination as recombination was detected when using the entire alignment even with only recombination free genes (Table 8). This finding may be explained by the recombination of fragments encompassing several genes or multiple recombination events involving smaller tracts along the genome. To test this hypothesis we used a method recently developed for the analysis of *Streptococcus agalactiae* genomes [75]. In order to identify patterns of recombination, nucleotide substitutions between strains were counted in sliding windows across the previously defined core chromosome representing 15 possible pair wise comparisons. Each pair wise comparison revealed highly conserved regions (<0.05% polymorphism on average) and less-conserved regions (>0.7% polymorphism), suggesting the occurrence of recombinational exchanges. When analyzing the different strains in depth we identified in each genome several regions with very low polymorphisms (below 0.05%) suggesting that DNA exchange of these fragments has occurred between the different *L. pneumophila* strains. Most interestingly, the two French strains Paris and HL 0604 1035 that are present since several years in France

Table 8 Intragenic and intergenic recombination in six *L. pneumophila* genomes predicted on individual genes and on combined data using six different methods

Data set	Event Number	Putative recombinant sequences	Detection Method					
			RDP	GENECONV	Boot scan	Max chi	Chimaera	SiScan
<i>metG</i>	1	Lorraine, Lens	NS	NS	NS	Yes	Yes	Yes
<i>dnaX</i>	1	Philadelphia	NS	NS	Yes	Yes	Yes	Yes
	2	Lens, Lorraine	NS	NS	NS	Yes	Yes	Yes
<i>proS</i>	1	HL06041035	Yes	Yes	Yes	Yes	Yes	Yes
	2	Philadelphia	NS	Yes	NS	Yes	Yes	Yes
<i>cysS</i>	1	Philadelphia	NS	NS	NS	Yes	Yes	NS
<i>lig</i>	1	Lorraine	NS	Yes	Yes	NS	NS	NS
<i>uvrC</i>	1	Lens, Philadelphia, Lorraine	NS	NS	NS	Yes	Yes	Yes
<i>flh</i>	1	Lens	NS	NS	Yes	Yes	Yes	Yes
	2	Paris, HL06041035	NS	NS	NS	Yes	Yes	Yes
<i>pgk</i>	1	Lens	NS	NS	NS	Yes	Yes	Yes
<i>atpD</i>	1	Corby	NS	NS	NS	Yes	Yes	Yes
Concatenated	1	Philadelphia	Yes	Yes	Yes	Yes	Yes	Yes
	2	Philadelphia	Yes	Yes	Yes	Yes	Yes	Yes
	3	HL06041035	Yes	Yes	Yes	Yes	Yes	Yes
	4	HL06041035	Yes	Yes	Yes	Yes	Yes	Yes
	5	Philadelphia, Corby, Lorraine	Yes	Yes	Yes	Yes	Yes	Yes
	6	Lens	Yes	Yes	Yes	Yes	Yes	Yes
	7	Paris, HL06041035	Yes	NS	NS	Yes	Yes	NS
	8	Paris	Yes	Yes	Yes	Yes	Yes	Yes
	9	Lens	Yes	Yes	NS	Yes	Yes	Yes
	10	Lens	Yes	Yes	Yes	Yes	Yes	NS
	11	HL06041035	Yes	Yes	NS	Yes	NS	NS
	12	Paris, HL06041035	Yes	Yes	Yes	Yes	Yes	Yes
	13	HL06041035, Lens	NS	Yes	NS	Yes	Yes	Yes
	14	Lens, Lorraine	Yes	NS	NS	Yes	Yes	NS
	15	Paris, HL06041035	Yes	Yes	NS	Yes	Yes	NS
	16	Corby	Yes	NS	NS	Yes	Yes	NS
	17	Lens	NS	Yes	NS	Yes	Yes	NS
	18	HL06041035, Paris	Yes	NS	NS	Yes	Yes	Yes
	19	Corby	Yes	NS	NS	Yes	Yes	NS
	20	Lorraine	Yes	Yes	NS	NS	NS	Yes
	21	Lens	Yes	NS	Yes	NS	NS	Yes
	22	Corby	Yes	NS	Yes	Yes	NS	Yes
	23	Lens	NS	Yes	NS	Yes	NS	NS
	24	Lens	NS	Yes	NS	Yes	NS	Yes
	25	Philadelphia	Yes	NS	NS	Yes	Yes	Yes

NS = non significant result. Yes = significant result with p-value ≤ 0.05 (where P is the highest acceptable probability value of recombination occurrence).

show 15 regions of a size between 10 and 99 kbs that have very low polymorphism and thus seem to have been exchanged between them (Additional file 7, Figure S1). In contrast when comparing strain Lens with the other 5 genomes analyzed here, very few regions with low polymorphism, two with strain HL 0604 1035 and one with strain Lorraine, were detected. Furthermore, no DNA exchanges seem to have occurred with strains Corby, Philadelphia or Paris. This indicates that strains

that are frequent in the same environment (e.g. strain Paris and HL 0604 1035) show high rates of DNA exchange probably by conjugation as suggested for *Streptococcus agalactiae* [75] and *Enterococcus faecalis* [76]. In contrast strain Lens, which has been identified to date only twice, in Lens (France) and in Germany, very few DNA transfers with the studied *L. pneumophila* strains seem to have taken place. Furthermore, some regions may be transferred also between several strains.

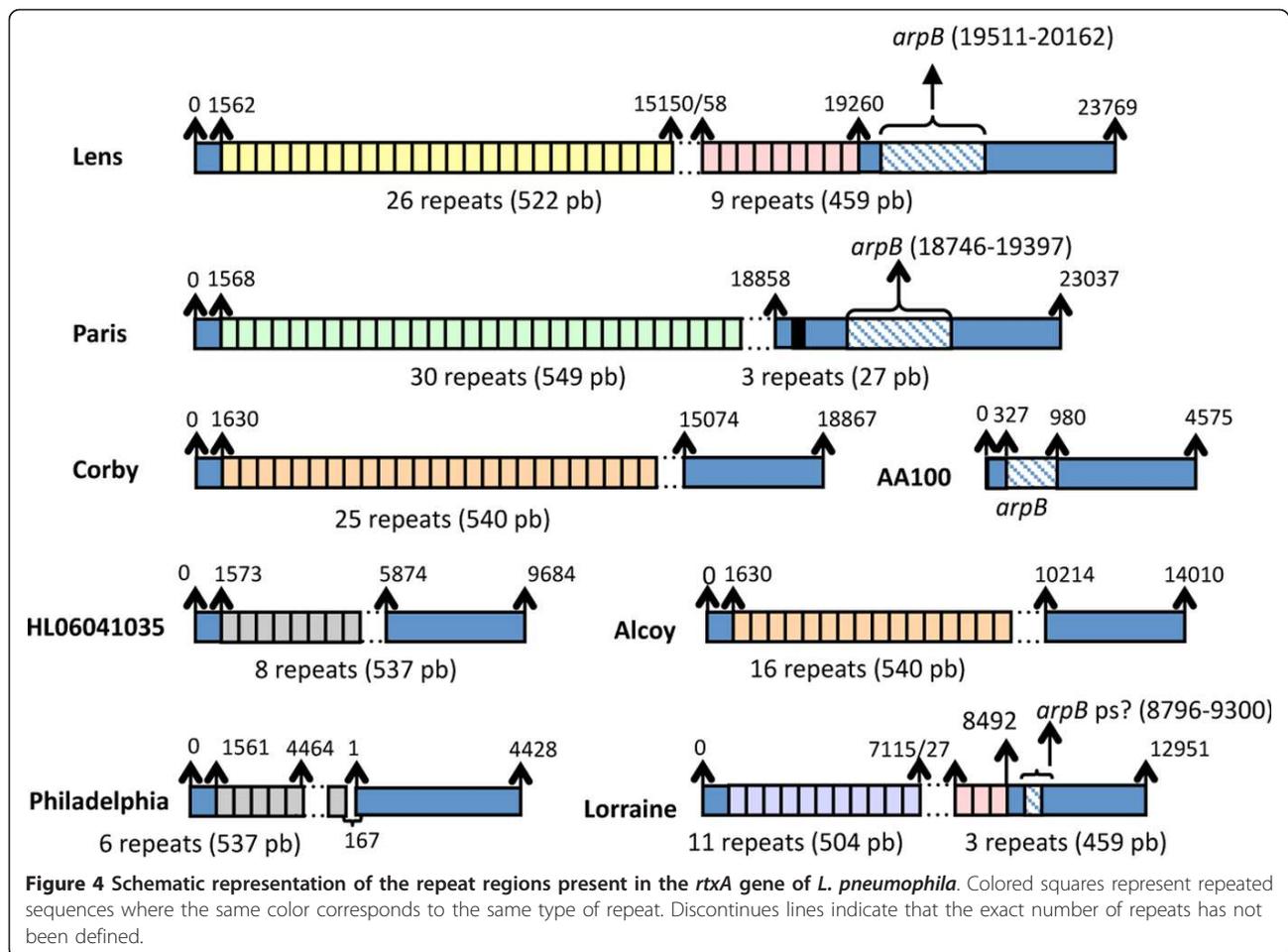
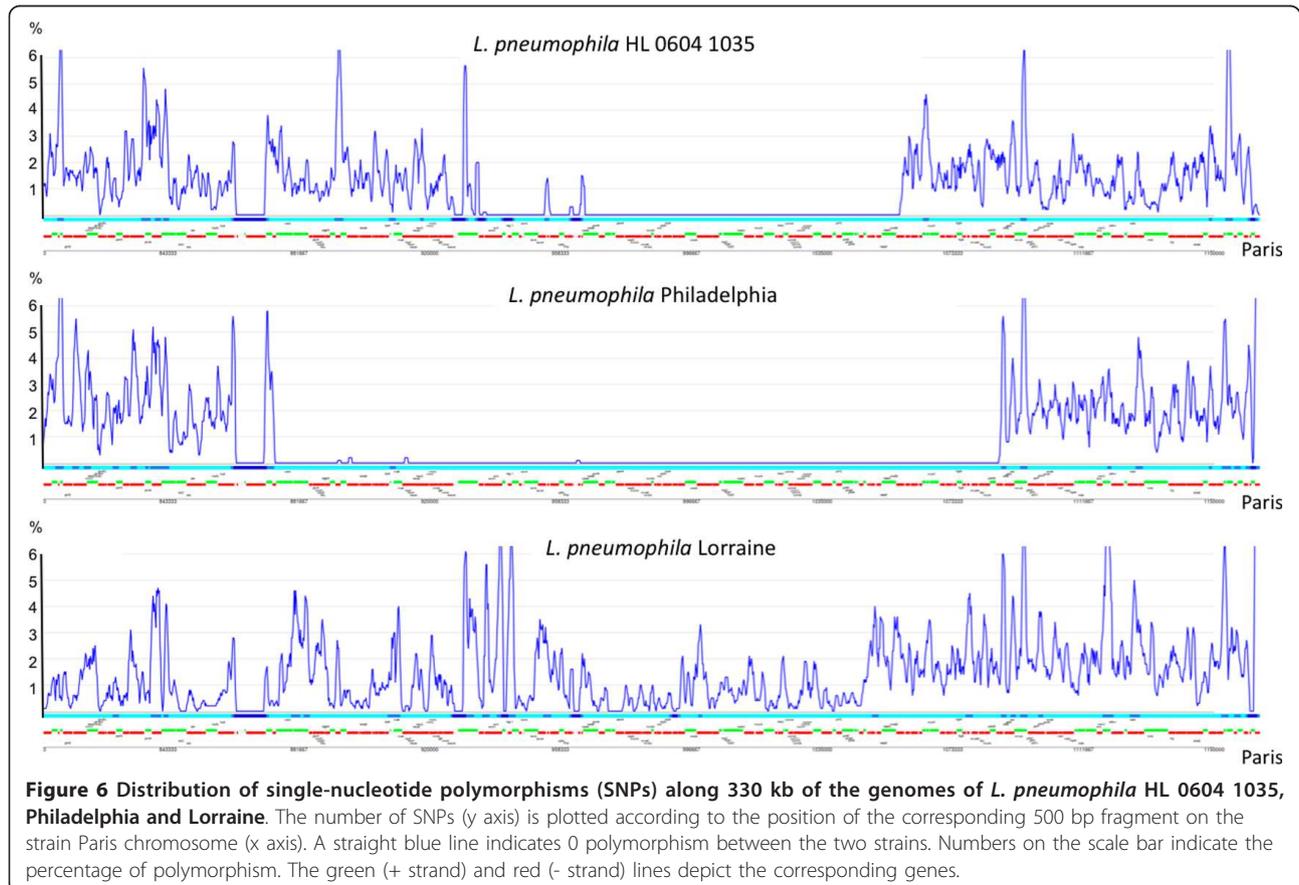
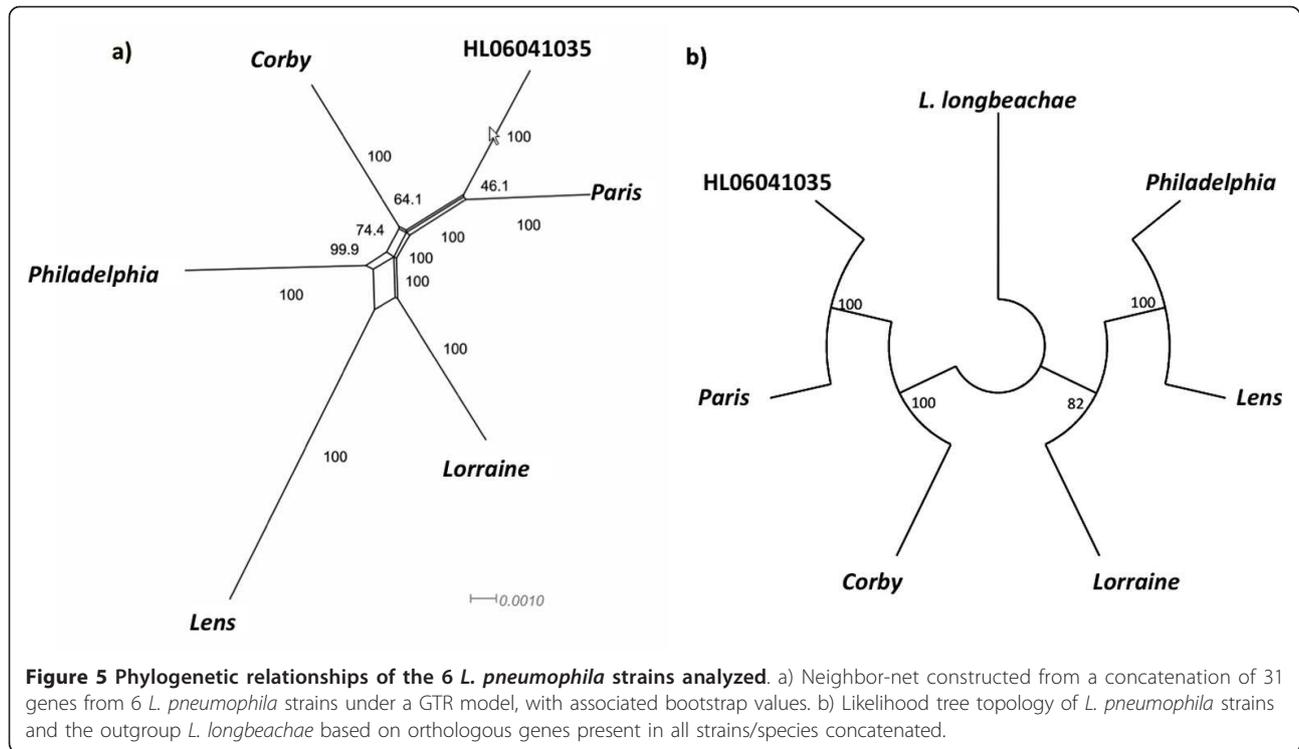


Figure 6 shows the distribution of single-nucleotide polymorphisms (SNPs) along 330 kb of the genome of *L. pneumophila* HL 0604 1035, Philadelphia and Lorraine as compared to the same region in the genome of strain Paris. We identified a region of 213 kbs a SNP frequency of 0.005%. Except an indel of 158 bs that shows higher polymorphism, only 11 SNPs are present in this region. This fragment may have evolved by conjugative transfer and recombination between strains Philadelphia and Paris. Among others, this region carries the genes necessary for lipopolysaccharide biosynthesis, that are also part of the smaller fragment that has been exchanged with strain HL 0604 1035. Our analyses suggest, that in addition to frequent intragenic recombination also recombination and horizontal transfer of large chromosomal fragments is taking place and shapes the chromosomes of *L. pneumophila*.

Conclusion

Analysis of the genome sequences of six *L. pneumophila* strains shows that the genomes of this environmental pathogen evolve by frequent HGT and high

recombination rates. Most interestingly, these events take place between eukaryotes and prokaryotes and among different strains and species of *Legionella*. A genome-wide map analysis of nucleotide polymorphisms among these six strains demonstrated that each chromosome is a mosaic of large chromosomal fragments from different origins suggesting that exchanges of large DNA regions of over 200 kb have contributed to the genome dynamics in the natural population. The many T4SS might be implicated in exchange of these fragments by conjugal transfer. Plasmids also play a role in genome diversification and are exchanged among strains and circulate even between different species of *Legionella*. Importantly, plasmids seem to excise and integrate into the genome probably depending on environmental cues. However, *L. pneumophila* encodes also several toxin anti-toxin that might help to stabilize certain mobile genetic elements. In the near future, the analyses of 100 s of genomes thanks to new generation sequencing combined with molecular studies should provide further clues about the genetic mechanisms and the evolutionary forces that shape the *Legionella* genomes.



Methods

Bacterial strains and sequence accession numbers

The strains sequenced in this study are *L. pneumophila* strain Lorraine [EMBL: FQ958210, EMBL:FQ958212] and *L. pneumophila* HL 0604 1035 [EMBL:FQ958211]. Strain Lorraine was isolated in 2004 from a patient and was recently described as a newly emerging endemic clone [26]. *L. pneumophila* strain HL 0604 1035 (ST 734, Bellingham subgroup of the Dresden panel) was isolated in 2006 from a water supply system in a French hospital that it is colonizing since more than 10 years.

Sequencing and assembly

The complete genome sequence of *L. pneumophila subsp. pneumophila* strain HL06041035 (A) and strain Lorraine (B) were determined using a Sanger/pyrosequencing hybrid approach. A shotgun library was constructed with 10kb size fragments, obtained after mechanical shearing of the total genomic DNA, and cloned into vector pCNS (pSU derived). Sequencing with vector-based primers was carried out using the ABI 3730 Applied Sequencer. A total of 20736 (A) and 21888 (B) reads (~4 fold-coverage) were analyzed and assembled with 502731 (A) and 555541 (B) reads (~15 fold-coverage) obtained with Genome Sequencer GS20 (Roche Applied Science). For the assembly, we used the Arachne "HybridAssemble" version (Broad Institute, <http://www.broad.mit.edu>) that combines the contigs obtained with 454 sequencing with Sanger reads. To validate the assembly, the Mekano interface (Genoscope), based on visualization of clone links inside and between contigs, was used to check the clone coverage and misassemblies. In addition, the consensus was confirmed using Consed functionalities <http://www.phrap.org>; the consensus quality and the high quality discrepancies. The finishing step was achieved by PCR, primer walking and *in vitro* transposition technology (Template Generation System™ II Kit; Finnzyme, Espoo, Finland), and a total of 930 (A) and 999 (B) sequences (109, 165 and 656 respectively for *L. pneumophila subsp. pneumophila* strain HL06041035 and 62, 204 and 733 respectively for *L. pneumophila subsp. pneumophila* str. Lorraine) were needed for gap closure and quality assessment.

Sequence analysis and annotation

The two newly sequenced *L. pneumophila* genomes were integrated into the MicroScope platform [77] to perform automatic and expert annotation of the genes, and comparative analysis with the other *L. pneumophila* strains already published. In addition the annotations of the previously published genomes were updated. The system integrates, for each predicted gene, the results of multiple bioinformatics methods (Blast result on

UniProt and specialized genomic data, InterPro, COG, PRIAM, synteny group computation using the complete bacterial genomes available at NCBI RefSeq, etc; more information on the syntactic and functional annotation process is given in [78]). In addition, many genomic and metabolic comparative tools are also available [77]. For details see <https://www.genoscope.cns.fr/agc/microscope/home/index.php>.

Definition of orthologous genes

To define orthologous chromosomal genes among the different *L. pneumophila* strains, pseudogenes and mobile elements were not taken into account due to the difficulty of ortholog assignment for these genes. Putative orthologous relations were defined as gene couples fulfilling two criteria: (i) having a bidirectional best hit (BBH) with an alignment threshold of 55% identity over at least 60% of the query sequence and target size (ii) and being in synteny. Subsequently, putative genes without any orthologous relation due to reduced identity percentage were integrated in a pre-existing orthologue group if they were flanked by orthologous genes showing gene order conservation (microsynteny). A final step of manual curation was carried out for each doubtful case.

Sequence alignments

For each gene of the selected data set, the nucleotide sequence was aligned based on the amino acid sequence using *tranalign*/EMBOSS package <http://emboss.sourceforge.net/>. Subsequently genes were concatenated in different data sets.

Identification of eukaryotic like proteins and eukaryotic domain carrying proteins

Eukaryotic domains were identified by analyzing the results obtained for all genes using the Interpro database that is integrated in MAGE. For the identification of eukaryotic like proteins we developed a new method. First we constructed two databases, one containing all and only eukaryotic sequences retrieved from public databases and a second one containing all and only prokaryotic sequences. From the second database we excluded the proteins of bacterial genera for which eukaryotic like protein-domains have been found in high proportions (e.g. parasites of protozoa) or bacterial genera that are reported to establish a symbiotic relationship with amoeba (for a detailed list see Additional file 8, Table S7). Those proteins, that showed a better, normalized blast score against eukaryotic proteins than to those present in the prokaryotic database were retrieved as eukaryotic like proteins. Parameters established for blast were: minimum identity: 25%; minimum

ratio avec query: 60%; minimum ratio avec target: 50%. The final results were manually checked.

Phylogenetic Analysis

For phylogenetic reconstruction of the *L. pneumophila* strains analyzed in this work several data sets were used: (i) 31 housekeeping genes described to be essential for all prokaryotes were selected based on the study of Zeigler [67] (Table 7 and Additional file 9, Figure S2) for a multi locus sequence analysis (MLSA) approach for which gene each was analyzed individually and as a concatenated alignment, (ii) a concatenated alignment of 2434 orthologous genes present in all analyzed *L. pneumophila* strains (iii) a concatenated alignment of 1867 orthologous genes present in all analyzed *L. pneumophila* strains and in the selected out group, *Legionella longbeachae* strain NSW150. An analysis of genetic divergence was performed using DNAsp vs 5.00.07 [79] using the 31 selected housekeeping genes. For phylogenetic reconstruction maximum likelihood (ML) methods were used to infer phylogenetic relationships for all data sets. Prior to ML analyses, a DNA substitution model for each gene or data set was selected using Modeltest v3.06 [80] and the Akaike information criterion. ML heuristic searches were performed using 500 random taxon-addition replicates with tree bisection and reconnection (TBR) and branch swapping. ML bootstrap support was determined using 1000 bootstrap replicates. The ML best trees were rooted on *L. longbeachae* when added. A network reconstruction was done for the same data set (i) using SplitsTree4 (version 4.10) [81]. The NeighborNet method and the GTR distance model were used to create the network.

Congruence test

The 31 genes selected for a MLST approach were tested for the significance of topological differences in the obtained phylogenetic trees using several methods. The first approach was based on the consensus of individual gene trees. The consensus tree was inferred using the CONSENSE program in the PHYLIP package <http://evolution.genetics.washington.edu/phytip.html> applying the extended majority rule. Secondly we tested the significance of topological differences in phylogenetic trees using the Shimodaira-Hasegawa (SH) test. The SH test compares the likelihood score (-lnL) of a given data set across its ML tree versus the -lnL of that data set across alternative topologies, which in this case are the ML phylogenies for other data sets. The differences in the -lnL values are evaluated for statistical significance using 1000 replicates based on resampling estimated with the log-likelihood (RELL) method (PAUP version 4.0b10; <http://paup.csit.fsu.edu/>). We applied the test using all

the trees obtained with individual genes, with the concatenated alignment against the alignment of each individual gene and with the alignment of all the 31 genes concatenated.

Recombination analysis

The 31 genes selected for a MLST approach and its corresponding concatenated alignment, were screened for the presence of putative recombination events by using RDP 2.0b08 [82]. This program identifies recombinant sequences and recombination breakpoints applying several methods. We selected six of them; two phylogenetic methods (which infer recombination when different parts of the genome result in discordant topologies): RDP [68], 2000) and Bootscanning [83]; and four nucleotide substitution methods (which examine the sequences either for a significant clustering of substitutions or for a fit to an expected statistical distribution): Maxchi and Chimaera [84], GeneConv [85] and Sis-scan [86]. We considered only those recombination events in our analysis that were identified by at least two methods. The common settings for all methods were (i) to consider sequences as circular, (ii) a statistical significance of $P < 0.05$, and (iii) a Bonferroni correction for multiple comparisons implemented in RDP.

Additional material

Additional file 1: Table S1: Nucleotide identity of 140 selected Dot/Icm substrates of strain Philadelphia and of their orthologs in the *L. pneumophila* strains analyzed in this study.

Additional file 2: Table S2: Genes specific of strain HL 0604 1035 with respect to strains Paris, Lens, Philadelphia, Corby and Lorraine.

Additional file 3: Table S3: Genes specific of strain Lorraine with respect to strains Paris, Lens, Philadelphia, Corby and HL0604 1035.

Additional file 4: Table S4: Summary of genetic diversity parameters for the 31 selected *L. pneumophila* genes used to establish the phylogeny.

Additional file 5: Table S5: Results for the SH Test of alternative topologies for the 6 analyzed *L. pneumophila* strains.

Additional file 6: Table S6: Conserved domains and repeats of the *rtxA* gene in 8 *L. pneumophila* strains.

Additional file 7: Figure S1 - Distribution of single-nucleotide polymorphisms (SNPs) along the genome of *L. pneumophila* HL 0604 1035 as compared to strains Lens, Philadelphia, Corby and Lorraine. The number of SNPs (y axis) is plotted according to the position of the corresponding 500 bp fragment on the strain Paris chromosome (x axis). A straight blue line indicates 0 polymorphism between the two strains. Numbers on the scale bar indicate the percentage of polymorphism. Yellow blocks indicate chromosomal regions with a SNP number lower than 0,005%.

Additional file 8: Tables S7 - List of bacterial genera removed from our prokaryotic database.

Additional file 9: Figure S2: Distribution of the 31 genes selected for establishing the phylogeny of *L. pneumophila* species. The coordinates are given with respect to the chromosome of *L. pneumophila* strain Paris. Numbers next to gene names indicate the first

position of the corresponding gene starting from the origin of replication.

Abbreviations

ANK: ankyrin motif; CRISPR: Clustered regularly interspaced short palindromic repeats; HGT: horizontal gene transfer; ML: maximum likelihood; nt: nucleotide; Sg1: serogroup 1; T4SS: Type IV secretion system; T2SS: Type II secretion system;

Acknowledgements

This work received financial support from the Institut Pasteur, the Centre National de la Recherche (CNRS), the Institut Carnot-Pasteur MI and from the ANR-10-PATH-004 project, in the frame of ERA-Net PathoGenoMics. The MicroScope platform got financial support from GIS IBIISA. L. Gomez-Valero was holder of a Roux postdoctoral research Fellowship financed by the Institut Pasteur and subsequently with support from the Fondation pour la Recherche Médicale (FRM). We would like to particularly thank Philippe Glaser for stimulating discussions and critical commenting of the article and Cyril Firmo for helping with the recombination analysis.

Author details

¹Institut Pasteur, Biologie des Bactéries Intracellulaires, 75724, Paris, France. ²CNRS URA 2171, 75724, Paris, France. ³Université de Lyon, Lyon, France, Centre National de Référence des Legionella, Lyon, France. ⁴INSERM, U851, 69007 Lyon, France. ⁵Hospices Civils de Lyon, Lyon, France. ⁶CEA/DSV/FAR/IG/Genoscope Laboratoire de Génomique Comparative, Evry Cedex, France. ⁷CNRS UMR8030 Laboratoire d'Analyses Bioinformatiques en Métabolisme et Génomique, Evry, France.

Authors' contributions

LGV and CB designed the study. SJ and JE supplied material and expertise; VB and BV performed genome sequencing; LGV and CR performed the genome annotation and analysis work, CM and RZ set up the LegioScope database. LGV and CB drafted and wrote the manuscript. All authors contributed to and approved the final manuscript.

Competing interests

The authors declare that they have no competing interests.

Received: 29 July 2011 Accepted: 1 November 2011

Published: 1 November 2011

References

1. Fliermans CB, Cherry WB, Orrison LH, Smith SJ, Tison DL, Pope DH: Ecological distribution of *Legionella pneumophila*. *Appl Environ Microbiol* 1981, **41**(1):9-16.
2. Fields BS: The molecular ecology of *Legionellae*. *Trends Microbiol* 1996, **4**(7):286-290.
3. Rowbotham TJ: Preliminary report on the pathogenicity of *Legionella pneumophila* for freshwater and soil amoebae. *J Clin Pathol* 1980, **33**(12):1179-1183.
4. Cazalet C, Rusniok C, Bruggemann H, Zidane N, Magnier A, Ma L, Tichit M, Jarraud S, Bouchier C, Vandenesch F, Kunst F, Etienne J, Glaser P, Buchrieser C: Evidence in the *Legionella pneumophila* genome for exploitation of host cell functions and high genome plasticity. *Nat Genet* 2004, **36**(11):1165-1173.
5. Hubber A, Roy CR: Modulation of host cell function by *Legionella pneumophila* type IV effectors. *Annu Rev Cell Dev Biol* 2010, **26**:261-283.
6. Nora T, Lomma M, Gomez-Valero L, Buchrieser C: Molecular mimicry: an important virulence strategy employed by *Legionella pneumophila* to subvert host functions. *Future Microbiol* 2009, **4**:691-701.
7. de Felipe KS, Pampou S, Jovanovic OS, Pericone CD, Ye SF, Kalachikov S, Shuman HA: Evidence for acquisition of *Legionella* type IV secretion substrates via interdomain horizontal gene transfer. *J Bacteriol* 2005, **187**(22):7716-7726.
8. Lurie-Weinberger MN, Gomez-Valero L, Merault N, Glockner G, Buchrieser C, Gophna U: The origins of eukaryotic-like proteins in *Legionella pneumophila*. *Int J Med Microbiol* 2010, **300**(7):470-481.
9. Brassinga AK, Hiltz MF, Sisson GR, Morash MG, Hill N, Garduno E, Edelstein PH, Garduno RA, Hoffman PS: A 65-Kilobase Pathogenicity Island Is Unique to Philadelphia-1 Strains of *Legionella pneumophila*. *J Bacteriol* 2003, **185**(15):4630-4637.
10. Cazalet C, Jarraud S, Ghavi-Helm Y, Kunst F, Glaser P, Etienne J, Buchrieser C: Multigenome analysis identifies a worldwide distributed epidemic *Legionella pneumophila* clone that emerged within a highly diverse species. *Genome Res* 2008, **18**(3):431-441.
11. Sexton JA, Vogel JP: Regulation of hypercompetence in *Legionella pneumophila*. *J Bacteriol* 2004, **186**(12):3814-3825.
12. Stone BJ, Kwaik YA: Natural competence for DNA transformation by *Legionella pneumophila* and its association with expression of type IV pili. *J Bacteriol* 1999, **181**(5):1395-1402.
13. Selander RK, McKinney RM, Whittam TS, Bibb WF, Brenner DJ, Nolte FS, Pattison PE: Genetic structure of populations of *Legionella pneumophila*. *J Bacteriol* 1985, **163**(3):1021-1037.
14. Edwards MT, Fry NK, Harrison TG: Clonal population structure of *Legionella pneumophila* inferred from allelic profiling. *Microbiology* 2008, **154**(Pt 3):852-864.
15. Coscolla M, Gosalbes MJ, Catalan V, Gonzalez-Candelas F: Genetic variability in environmental isolates of *Legionella pneumophila* from Comunidad Valenciana (Spain). *Environ Microbiol* 2006, **8**(6):1056-1063.
16. Bumbaugh AC, McGraw EA, Page KL, Selander RK, Whittam TS: Sequence polymorphism of *dotA* and *mip* alleles mediating invasion and intracellular replication of *Legionella pneumophila*. *Curr Microbiol* 2002, **44**(5):314-322.
17. Ko KS, Hong SK, Lee HK, Park MY, Kook YH: Molecular evolution of the *dotA* gene in *Legionella pneumophila*. *J Bacteriol* 2003, **185**(21):6269-6277.
18. Ko KS, Lee HK, Park MY, Park MS, Lee KH, Woo SY, Yun YJ, Kook YH: Population genetic structure of *Legionella pneumophila* inferred from RNA polymerase gene (*rpoB*) and *DotA* gene (*dotA*) sequences. *J Bacteriol* 2002, **184**(8):2123-2130.
19. Coscolla M, Gonzalez-Candelas F: Population structure and recombination in environmental isolates of *Legionella pneumophila*. *Environ Microbiol* 2007, **9**(3):643-656.
20. Vos M, Didelot X: A comparison of homologous recombination rates in bacteria and archaea. *ISME J* 2009, **3**(2):199-208.
21. Coscolla M, Comas I, Gonzalez-Candelas F: Quantifying nonvertical inheritance in the evolution of *Legionella pneumophila*. *Mol Biol Evol* 2011, **28**(2):985-1001.
22. Feil EJ, Spratt BG: Recombination and the population structures of bacterial pathogens. *Annu Rev Microbiol* 2001, **55**:561-590.
23. Harrison TG, Doshi N, Fry NK, Joseph CA: Comparison of clinical and environmental isolates of *Legionella pneumophila* obtained in the UK over 19 years. *Clin Microbiol Infect* 2007, **13**(1):78-85.
24. Steinert M, Heuner K, Buchrieser C, Albert-Weissenberger C, Glockner G: *Legionella* pathogenicity: genome structure, regulatory networks and the host cell response. *Int J Med Microbiol* 2007, **297**(7-8):577-587.
25. Chien M, Morozova I, Shi S, Sheng H, Chen J, Gomez SM, Asamani G, Hill K, Nuara J, Feder M, Rineer J, Greenberg JJ, Steshenko V, Park SH, Zhao B, Teplitskaya E, Edwards JR, Pampou S, Georgiou A, Chou IC, Iannuccilli W, Ulz ME, Kim DH, Geringer-Sameth A, Goldsberry C, Morozov P, Fischer SG, Segal G, Qu X, Rzhetsky A, et al: The genomic sequence of the accidental pathogen *Legionella pneumophila*. *Science* 2004, **305**(5692):1966-1968.
26. Ginevra C, Forey F, Campese C, Reyrolle M, Che D, Etienne J, Jarraud S: Lorraine strain of *Legionella pneumophila* serogroup 1, France. *Emerg Infect Dis* 2008, **14**(4):673-675.
27. Schmitz-Esser S, Tischler P, Arnold R, Montanaro J, Wagner M, Rattei T, Horn M: The genome of the amoeba symbiont "Candidatus *Amoebophilus asiaticus*" reveals common mechanisms for host cell interaction among amoeba-associated bacteria. *J Bacteriol* 2010, **192**(4):1045-1057.
28. Cazalet C, Gomez-Valero L, Rusniok C, Lomma M, Dervins-Ravault D, Newton HJ, Sansom FM, Jarraud S, Zidane N, Ma L, Bouchier C, Etienne J, Hartland EL, Buchrieser C: Analysis of the *Legionella longbeachae* genome and transcriptome uncovers unique strategies to cause Legionnaires' disease. *PLoS Genet* 2010, **6**(2):e1000851.
29. Hunter S, Apweiler R, Attwood TK, Bairoch A, Bateman A, Binns D, Bork P, Das U, Daugherty L, Duquenne L, Finn RD, Gough J, Haft D, Hulo N, Kahn D, Kelly E, Laugraud A, Letunic I, Lonsdale D, Lopez R, Madera M, Maslen J, McAnulla C, McDowall J, Mistry J, Mitchell A, Mulder N, Natale D,

- Orengo C, Quinn AF, et al: InterPro: the integrative protein signature database. *Nucleic Acids Res* 2009(D):211-215.
30. Djinovic-Carugo K, Gautel M, Ylänne J, Young P: The spectrin repeat: a structural platform for cytoskeletal protein assemblies. *FEBS Lett* 2002, **513**(1):119-123.
31. de Felipe KS, Glover RT, Charpentier X, Anderson OR, Reyes M, Pericone CD, Shuman HA: *Legionella* eukaryotic-like type IV substrates interfere with organelle trafficking. *PLoS Pathog* 2008, **4**(8):e1000117.
32. Heidtman M, Chen EJ, Moy MY, Isberg RR: Large-scale identification of *Legionella pneumophila* Dot/Icm substrates that modulate host cell vesicle trafficking pathways. *Cell Microbiol* 2009, **11**(2):230-248.
33. Isberg RR, O'Connor TJ, Heidtman M: The *Legionella pneumophila* replication vacuole: making a cosy niche inside host cells. *Nat Rev Microbiol* 2009, **7**(1):13-24.
34. Paduch M, Jeleń F, Otlewski J: Structure of small G proteins and their regulators. *Acta Biochim Pol* 2001, **48**(8):829-850.
35. Barr F, Lambright DG: Rab GEFs and GAPs. *Curr Opin Cell Biol* 2010, **22**(4):461-470.
36. Ninio S, Celli J, Roy CR: A *Legionella pneumophila* Effector Protein Encoded in a Region of Genomic Plasticity Binds to Dot/Icm-Modified Vacuoles. *PLoS Pathog* 2009, **5**(1):e1000278.
37. Berger KH, Isberg RR: Two distinct defects in intracellular growth complemented by a single genetic locus in *Legionella pneumophila*. *Mol Microbiol* 1993, **7**(1):7-19.
38. Segal G, Shuman HA: Characterization of a new region required for macrophage killing by *Legionella pneumophila*. *Infect Immun* 1997, **65**(12):5057-5066.
39. Morozova I, Qu X, Shi S, Asamani G, Greenberg JE, Shuman HA, Russo JJ: Comparative sequence analysis of the *icm/dot* genes in *Legionella*. *Plasmid* 2004, **51**(2):127-147.
40. Nagai H, Kagan JC, Zhu X, Kahn RA, Roy CR: A bacterial guanine nucleotide exchange factor activates ARF on *Legionella* phagosomes. *Science* 2002, **295**(5555):679-682.
41. Burstein D, Zusman T, Degtyar E, Viner R, Segal G, Pupko T: Genome-scale identification of *Legionella pneumophila* effectors using a machine learning approach. *PLoS Pathog* 2009, **5**(7):e1000508.
42. Campodonico EM, Chesnel L, Roy CR: A yeast genetic system for the identification and characterization of substrate proteins transferred into host cells by the *Legionella pneumophila* Dot/Icm system. *Mol Microbiol* 2005, **56**(4):918-933.
43. Shohdy N, Efe JA, Emr SD, Shuman HA: Pathogen effector protein screening in yeast identifies *Legionella* factors that interfere with membrane trafficking. *Proc Natl Acad Sci USA* 2005, **102**(13):4866-4871.
44. Zhu W, Banga S, Tan Y, Zheng C, Stephenson R, Gately J, Luo ZQ: Comprehensive Identification of Protein Substrates of the Dot/Icm Type IV Transporter of *Legionella pneumophila*. *PLoS One* 2011, **6**(3):e17638.
45. Ivanov SS, Roy CR: Modulation of ubiquitin dynamics and suppression of DALIS formation by the *Legionella pneumophila* Dot/Icm system. *Cell Microbiol* 2009, **11**(2):261-278.
46. Price CT, Al-Quadan T, Santic M, Jones SC, Abu Kwaik Y: Exploitation of conserved eukaryotic host cell farnesylation machinery by an F-box effector of *Legionella pneumophila*. *J Exp Med* 2010, **207**(8):1713-1726.
47. Juhas M, Crook DW, Hood DW: Type IV secretion systems: tools of bacterial horizontal gene transfer and virulence. *Cell Microbiol* 2008, **10**(12):2377-2386.
48. Christie PJ, Atmakuri K, Krishnamoorthy V, Jakubowski S, Cascales E: Biogenesis, architecture, and function of bacterial type IV secretion systems. *Annu Rev Microbiol* 2005, **59**(451-485).
49. Backert S, Meyer TF: Type IV secretion systems and their effectors in bacterial pathogenesis. *Curr Opin Microbiol* 2006, **9**(2):207-217.
50. Frost LS, Ippen-Ihler K, Skurray RA: Analysis of the sequence and gene products of the transfer region of the F sex factor. 1994.
51. Ogata H, La Scola B, Audic S, Renesto P, Blanc G, Robert C, Fournier PE, Claverie JM, Raoult D: Genome sequence of *Rickettsia bellii* illuminates the role of amoebae in gene exchanges between intracellular pathogens. *PLoS Genet* 2006, **2**(5):e76.
52. Doleans-Jordheim A, Akermi M, Ginevra C, Cazalet C, Kay E, Schneider D, Buchrieser C, Atlan D, Vandenesch F, Etienne J, Jarraud S: Growth-phase-dependent mobility of the *Ivh*-encoding region in *Legionella pneumophila* strain Paris. *Microbiology* 2006, **152**(Pt 12):3561-3568.
53. Lawley TD, Klimke WA, Gubbins MJ, Frost LS: F factor conjugation is a true type IV secretion system. *FEMS Microbiol Lett* 2003, **224**(1):1-15.
54. Glöckner G, Albert-Weissenberger C, Weinmann E, Jacobi S, Schunder E, Steinert M, Hacker J, Heuner K: Identification and characterization of a new conjugation/type IVA secretion system (*trb/tra*) of *Legionella pneumophila* Corby localized on a mobile genomic island. *Int J Med Microbiol* 2007, **298**(5-6):411-428.
55. Molofsky AB, Swanson MS: *Legionella pneumophila* CsrA is a pivotal repressor of transmission traits and activator of replication. *Mol Microbiol* 2003, **50**(2):445-461.
56. Sorek R, Kunin V, Hugenholz P: CRISPR—a widespread system that provides acquired resistance against phages in bacteria and archaea. *Nat Rev Microbiol* 2008, **6**(3):181-186.
57. Van Melderen L: Toxin-antitoxin systems: why so many, what for? *Curr Opin Microbiol* 2010, **13**(6):781-785.
58. Christensen-Dalsgaard M, Gerdes K: Two *higBA* loci in the *Vibrio cholerae* superintegron encode mRNA cleaving enzymes and can stabilize plasmids. *Mol Microbiol* 2006, **62**(2):397-411.
59. Szekeres S, Dauti M, Wilde C, Mazel D, Rowe-Magnus DA: Chromosomal toxin-antitoxin loci can diminish large-scale genome reductions in the absence of selection. *Mol Microbiol* 2007, **63**(6):1588-1605.
60. Jore MM, Brouns SJ, van der Oost J: RNA in Defense: CRISPRs Protect Prokaryotes against Mobile Genetic Elements. *Cold Spring Harb Perspect Biol* 2011.
61. D'Auria G, Jimenez-Hernandez N, Peris-Bondia F, Moya A, Latorre A: *Legionella pneumophila* pangenome reveals strain-specific virulence factors. *BMC Genomics* 2010, **11**(1):181.
62. Lammertyn E, Vande Voorde J, Meyen E, Maes L, Mast J, Anné J: Evidence for the presence of *Legionella* bacteriophages in environmental water samples. *Microb Ecol* 2008, **56**(1):191-197.
63. Pragai Z, Harwood CR: Regulatory interactions between the Pho and sigma(B)-dependent general stress regulons of *Bacillus subtilis*. *Microbiology* 2002, **148**(Pt 5):1593-1602.
64. Collingro A, Toenshoff ER, Taylor MW, Fritsche TR, Wagner M, Horn M: 'Candidatus *Protochlamydia amoebophila*', an endosymbiont of *Acanthamoeba* spp. *Int J Syst Evol Microbiol* 2005, **55**(Pt 5):1863-1866.
65. La Scola B, Birtles RJ, Greub G, Harrison TJ, Ratcliff RM, Raoult D: *Legionella drancourtii* sp. nov., a strictly intracellular amoebal pathogen. *Int J Syst Evol Microbiol* 2004, **54**(Pt 3):699-703.
66. Kim BR, Anderson JE, Mueller SA, Gaines WA, Kendall AM: Literature review—efficacy of various disinfectants against *Legionella* in water systems. *Water Res* 2002, **36**(18):4433-4444.
67. Zeigler DR: Gene sequences useful for predicting relatedness of whole genomes in bacteria. *Int J Syst Evol Microbiol* 2003, **53**(Pt 6):1893-1900.
68. Martin D, Rybicki E: RDP: detection of recombination amongst aligned sequences. *Bioinformatics* 2000, **16**(6):562-563.
69. Ko KS, Lee HK, Park MY, Lee KH, Yun YJ, Woo SY, Miyamoto H, Kook YH: Application of RNA polymerase beta-subunit gene (*rpoB*) sequences for the molecular differentiation of *Legionella* species. *J Clin Microbiol* 2002, **40**(7):2653-2658.
70. Coscolla M, Gonzalez-Candelas F: Comparison of clinical and environmental samples of *Legionella pneumophila* at the nucleotide sequence level. *Infect Genet Evol* 2009, **9**(5):882-888.
71. Ko KS, Lee HK, Park MY, Kook YH: Mosaic structure of pathogenicity islands in *Legionella pneumophila*. *J Mol Evol* 2003, **57**(1):63-72.
72. Ko KS, Miyamoto H, Lee HK, Park MY, Fukuda K, Park BJ, Kook YH: Genetic diversity of *Legionella pneumophila* inferred from *rpoB* and *dotA* sequences. *Clin Microbiol Infect* 2006, **12**(3):254-261.
73. D'Auria G, Jiménez N, Peris-Bondia F, Pelaz C, Latorre A, Moya A: Virulence factor *rtx* in *Legionella pneumophila*, evidence suggesting it is a modular multifunctional protein. *BMC Genomics* 2008, **9**:14.
74. Posada D, Crandall KA: The effect of recombination on the accuracy of phylogeny estimation. *J Mol Evol* 2002, **54**(3):396-402.
75. Brochet M, Rusniok C, Couve E, Dramsi S, Poyart C, Trieu-Cuot P, Kunst F, Glaser P: Shaping a bacterial genome by large chromosomal rearrangements, the evolutionary history of *Streptococcus agalactiae*. *Proc Natl Acad Sci USA* 2008, **105**(41):15961-15966.
76. Manson JM, Hancock LE, Gilmore MS: Mechanism of chromosomal transfer of *Enterococcus faecalis* pathogenicity island, capsule, antimicrobial resistance, and other traits. *Proc Natl Acad Sci USA* 107(27):12269-12274.

77. Vallenet D, Labarre L, Rouy Z, Barbe V, Bocs S, Cruveiller S, Lajus A, Pascal G, Scarpelli C, Medigue C: **MaGe: a microbial genome annotation system supported by synteny results.** *Nucleic Acids Res* 2006, **34**(1):53-65.
78. Vallenet D, Engelen S, Mornico D, Cruveiller S, Fleury L, Lajus A, Rouy Z, Roche D, Salvignol G, Scarpelli C, Médigue C: **MicroScope: a platform for microbial genome annotation and comparative genomics.** *Database (Oxford)* 2009, bap021.
79. Rozas J, Sánchez-DelBarrio JC, Messeguer X, Rozas R: **DnaSP, DNA polymorphism analyses by the coalescent and other methods.** *Bioinformatics* 2003, **19**:18.
80. Posada D, Crandall KA: **MODELTEST: testing the model of DNA substitution.** *Bioinformatics* 1998, **14**(9):817-818.
81. Huson DH, Bryant D: **Application of phylogenetic networks in evolutionary studies.** *Mol Biol Evol* 2006, **23**(2):254-267.
82. Martin DP, Williamson C, Posada D: **RDP2: recombination detection and analysis from sequence alignments.** *Bioinformatics* 2005, **21**(2):260-262.
83. Salminen MO, Carr JK, Burke DS, McCutchan FE: **Identification of breakpoints in intergenotypic recombinants of HIV type 1 by bootscanning.** *AIDS Res Hum Retroviruses* 1995, **11**(11):1423-1425.
84. Smith JM: **Analyzing the mosaic structure of genes.** *J Mol Evol* 1992, **34**(2):126-129.
85. Padidam M, Beachy RN, Fauquet CM: **A phage single-stranded DNA (ssDNA) binding protein complements ssDNA accumulation of a geminivirus and interferes with viral movement.** *J Virol* 1999, **73**(2):1609-1616.
86. Gibbs MJ, Armstrong JS, Gibbs AJ: **Sister-scanning: a Monte Carlo procedure for assessing signals in recombinant sequences.** *Bioinformatics* 2000, **16**(7):573-582.

doi:10.1186/1471-2164-12-536

Cite this article as: Gomez-Valero et al.: Extensive recombination events and horizontal gene transfer shaped the *Legionella pneumophila* genomes. *BMC Genomics* 2011 **12**:536.

**Submit your next manuscript to BioMed Central
and take full advantage of:**

- Convenient online submission
- Thorough peer review
- No space constraints or color figure charges
- Immediate publication on acceptance
- Inclusion in PubMed, CAS, Scopus and Google Scholar
- Research which is freely available for redistribution

Submit your manuscript at
www.biomedcentral.com/submit

